

UNCLASSIFIED

AD NUMBER

AD467808

LIMITATION CHANGES

TO:

Approved for public release; distribution is unlimited.

FROM:

Distribution authorized to U.S. Gov't. agencies and their contractors;
Administrative/Operational Use; JUN 1965. Other requests shall be referred to Defense Advanced Research Projects Agency, 675 North Randolph Street, Arlington, VA 22203-2114.

AUTHORITY

ARPA per RAND ltr dtd 26 Aug 1970

THIS PAGE IS UNCLASSIFIED

SECURITY

MARKING

The classified or limited status of this report applies to each page, unless otherwise marked.

Separate page printouts MUST be marked accordingly.

THIS DOCUMENT CONTAINS INFORMATION AFFECTING THE NATIONAL DEFENSE OF THE UNITED STATES WITHIN THE MEANING OF THE ESPIONAGE LAWS, TITLE 18, U.S.C., SECTIONS 793 AND 794. THE TRANSMISSION OR THE REVELATION OF ITS CONTENTS IN ANY MANNER TO AN UNAUTHORIZED PERSON IS PROHIBITED BY LAW.

NOTICE: When government or other drawings, specifications or other data are used for any purpose other than in connection with a definitely related government procurement operation, the U. S. Government thereby incurs no responsibility, nor any obligation whatsoever; and the fact that the Government may have formulated, furnished, or in any way supplied the said drawings, specifications, or other data is not to be regarded by implication or otherwise as in any manner licensing the holder or any other person or corporation, or conveying any rights or permission to manufacture, use or sell any patented invention that may in any way be related thereto.

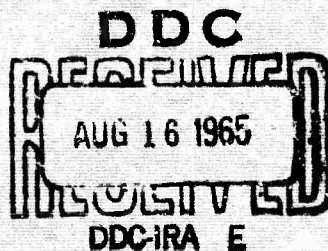
MEMORANDUM

RM-4511-ARPA

JUNE 1965

THE RELIABILITY OF
GROUND-BASED DIGITAL COMPUTERS

Rodger R. Lowe and Michael Warshaw



PREPARED FOR:

ADVANCED RESEARCH PROJECTS AGENCY

The RAND Corporation
SANTA MONICA • CALIFORNIA

MEMORANDUM
RM-4511-ARPA
JUNE 1965

THE RELIABILITY OF GROUND-BASED DIGITAL COMPUTERS

Rodger R. Lowe and Michael Warshaw

This research is supported by the Advanced Research Projects Agency under Contract No. SD-79. Any views or conclusions contained in this Memorandum should not be interpreted as representing the official opinion or policy of ARPA.

DDC AVAILABILITY NOTICE

Qualified requesters may obtain copies of this report from the Defense Documentation Center (DDC).

The RAND Corporation

1700 MAIN ST. • SANTA MONICA • CALIFORNIA • 90406

NOT approved for ODS release

PREFACE

This Memorandum was prepared for the Advanced Research Projects Agency (ARPA) as part of RAND's continued interest in the problems of data processing for ballistic missile defense.

The original intent was to write specifically on the topic of computer reliability for missile defense, but the scope increased to cover the broader field of the operational availability of data processing systems, and in particular, those ground-based systems which can be repaired.

Selected portions of this report should be of special interest to statisticians, circuit designers, and programmers; and the conclusions should interest operations and systems analysts and others who are concerned with computer reliability.

This work should not be construed as a handbook on data processor reliability; the treatment of the topics is not uniform, the authors having been guided to a large extent by their own opinions and interests. We ask, therefore, the tolerance of any specialist whose particular area of interest may appear slighted. Space does not permit completeness in all the areas which contribute to a reliable computer.

One of the authors, Rodger Lowe, is Vice President of the Mesa Scientific Corporation of Inglewood, California and a consultant to The RAND Corporation.

SUMMARY

The electronic computer plays a steadily increasing role in the affairs of man; and man, in turn, relies more and more on the voice of the computer. As this trend continues, the computer grows in size and cost, and reliability becomes not just a desirable feature, but an absolute necessity. The effects of a computer failing during a time of national emergency or when the control and safety of an entire steel mill are at stake can be catastrophic.

This Memorandum discusses the many aspects (both qualitative and quantitative) of obtaining a reliable digital computer and, in particular, investigates that class of ground-based data processing systems where repair is possible.

The study begins by reviewing the reliability of computer parts (transistors, capacitors, integrated circuits) and applies the results to a large variety of probabilistic models of system availability. Further, it discusses the availability of ground-based data processing systems--specifically, the probability that a repairable computer which should be ready will in fact be ready for use at some arbitrary future time. It is concluded that part failure distributions show a form of decreasing failure rates for the entire population which in no way correlates with the predicted behavior of the ideal part. The total part population shows a decreasing failure rate because, and only because, various controllably small subgroups show initially increasing failure rates until every member

of the subgroup has failed, at which time the failure rate of the entire population effectively decreases.

Next, the authors survey the various machine structures which yield higher reliability and show that, where service is available, redundancy is never a contender as a means to high reliability. With service, the multi-processor provides the highest availability, the multiple-processor (duplex or triplex), second highest.

A prediction of the failure rates for the best parts available in 1968 yields values which are roughly about a factor of three smaller than the best 1965 values, and a factor of ten smaller than what are considered "good" 1965 parts. Significant improvements in machine availability will be realized as a result of this decrease in part failure rate.

Computer availability can also be improved by decreasing the time spent in repairing a faulty machine. A survey of developments in automatic fault diagnosis and isolation reveals that systems of considerable elegance and power are now available, though expensive. Their employment would considerably reduce service time.

Chapter VII gives a much fuller summary of this Memorandum. There, also, the reader will find references into the text for explicit recommendations on improving the reliability of the computer and its associated programs.

CONTENTS

PREFACE	iii
SUMMARY	v
TABLES	xi
FIGURES	xiii
SYMBOLS	xvii
Chapter	
I. INTRODUCTION	1
1. Definitions	3
2. The Scope of the Report	4
II. THE PROCESS OF FAILURE	8
1. Introduction	8
2. Failure Modes and Mechanisms	18
3. Models of Failure Behavior	25
4. Estimates of Non-Constant Failure Rates	32
5. Storage Life Versus Operating Life .	35
6. Cost-Reliability Relationships	37
7. Present and Predicted Reliability Levels	49
8. Predicted Part Costs	56
Improvements in Materials, Processes, and Quality Control .	57
Possible Breakthroughs	58
9. Conclusions and Recommendations on Parts	59
Semiconductor Devices	59
Resistors	62
Capacitors	62
In-Plant Parts Handling	63
III. THE RELIABLE COMPUTER	66
1. Introduction	66
2. Circuit Design	67
Bogey Design	69
Worst-Case Design	70

Statistical Design	71
Protection Against Part Failure ..	72
Recommendations in the Procurement of Good Circuit Design	75
3. Logical Design	76
4. Circuit Redundancy	79
5. Failure Detection	84
IV. MULTIPLE COMPUTERS FOR RELIABILITY	94
1. Introduction	94
2. The Multi-Processor	95
3. System Comparisons	96
V. SYSTEMS CONSIDERATIONS	102
1. Introduction	102
2. Programming	102
Techniques for Avoiding Errors ...	104
3. Programmed Error Detection	106
4. Interconnection and Packaging Reliability	107
5. Extreme Environments	112
6. The Role of the Manufacturer	113
VI. THE PROCESS OF MAINTENANCE	114
1. Introduction	114
2. Diagnostic Techniques	116
Method 1	118
Method 2	119
Method 3	121
3. The Maintenance Module	123
4. Preventive Maintenance	128
VII. SUMMARY	131
1. Introduction	131
2. Parts and a Definition of Failure ..	132
3. Part Failure Mechanisms and Distributions	133
4. Standby Conditions	134
5. The "Equivalent Transistor" Computer and Some Typical Systems	134
6. The Unit Parts Complement	136
7. Predicted Part Failure Rates	136
8. Integrated Circuits	136
9. Circuit Design	136

10.	Part Redundancy	138
11.	Failure Detection	138
12.	Multiple Computers and the Multi-Processor	138
13.	Programming	141
14.	Connectors and Packaging	146
15.	The Manufacturer	147
16.	Automatic Fault Diagnosis and Isolation	147
17.	Optimum Module Size	147

APPENDIX

A. INTRODUCTION TO THE MATHEMATICS OF

AVAILABILITY	149
1. Introduction	149
2. Definitions and the Poisson Process	152
3. No Redundancy--No Service and Exponential Failure Distribution .	157
4. $2m+1$ -Fold Redundancy with Perfect Voting--No Service and Exponential Failure Distribution	159
5. N-Fold Redundancy with Imperfect Voting--No Service and Exponential Failure Distribution	162
6. No Redundancy--Exponential Service and Failure Distribution	173
7. Redundant Computer--Exponential Service and Failure Distributions	175
8. Non-Redundant Computer--Exponential Service and Weibull Failure Distributions	197
9. Multiple Computers	200
10. Single Computer--Exponential Service and Failure Distributions	200
11. Duplex Computers with Flexible Service--Exponential Service and Failure Distributions	202
12. Duplex Computers with Limited Service--Exponential Service and Failure Distributions	203
13. Triplex Computers with Flexible Service--Exponential Service and Failure Distributions	209
14. The Multi-Processor System--Exponential Service and Failure Distributions	209

B.	SEMICONDUCTOR DEVICES--BEHAVIOR VS.	
	STRESS	225
1.	Introduction	225
2.	Factors Contributing to Semi-	
	conductor Failure	226
	Bulk Factors	226
	Contact Factors	227
	Packaging Factors	227
3.	Semiconductor Failures Listed by	
	Cause	227
	Gross Manufacturing, Design, and	
	Mechanical Failures	228
	Contact and Interconnection	
	Failures	228
	Failures Due to the Surface and	
	Surface Environment	229
	Bulk and Process Failures	229
	Gross Externally Induced	
	Failures	229
	Integrated Circuits and Tran-	
	sistors--Failure Origins, Modes,	
	and Mechanisms	230
	Semiconductor Device Degradation .	230
C.	RESISTORS--BEHAVIOR VS. STRESS	235
1.	Carbon Composition Resistors	235
2.	Metal and Carbon Film Resistors	237
3.	Tin Oxide Resistors	239
D.	CAPACITORS--BEHAVIOR VS. STRESS	241
1.	Dipped Mica Capacitors	241
2.	Glass Capacitors	244
3.	Paper and Electrolytic Capacitors ..	245
E.	A COMPENDIUM OF FAILURE STATISTICS	247
	REFERENCES	255

TABLES

Table

II-1	Examples of Integrated Circuit Failure	19
II-2	Normally Distributed Wearout	24
II-3	Failure Rate vs. Time for Synthetic Failure Distribution of Fig. II-1	29
II-4	Estimated Parameters for Weibull Failure Distribution	33
II-5	Costs for Three Grades of Parts	41
II-6	Complexity of Existing Systems	42
II-7	Part Percentages of an Average Computer	43
II-8	Cost vs. Procurement Policy for a Unit Parts Complement	44
II-9	Failure Rate vs. Procurement Policy For a Unit Parts Complement Using Discrete Parts	44
II-10	Distribution of Parts in a Unit System .	45
II-11	Failure Rate vs. Procurement Policy For a Unit Parts Complement Using Integrated Circuits	47
II-12	Summary of Cost-Reliability Results	48
II-13	Predicted Part Failure Rates (%/1000 hr, 10-year average)	50
II-14	Per Cent Contribution of Planar Process Defects to Part Failure	51
II-15	Relative Failure Rate for Integrated Circuits	53
II-16	Predicted Price of Integrated Circuits .	56
IV-1	Size-Reliability Factor for Various Part Grades	101
V-1	Connections for Etched Board and Integrated Circuits	109

VI-1	Spare Parts Cost vs. Number of Modules	125
VI-2	Estimated Diagnostic Time vs. Number of Modules	127
VII-1	Complexity of Existing Systems	135
VII-2	Predicted Part Failure Rates (%/1000 hr, 10-year average)	137
B-1	Integrated Circuits and Transistors Failure Origins, Modes, and Mechanisms	231
B-2	Failure Origin, Mode, and Mechanism Key	232
B-3	Life Test Data on h_{FE} and I_{CBO} Variation	233
C-1	Carbon Composition Resistors-- Behavior vs. Stress	236
C-2	Test Data on Type N20 Tin Oxide Resistor	239
C-3	Test Data on Type A-100 Tin Oxide Resistor	240
C-4	Reliability of Corning Tin Oxide Resistors	240
D-1	Temperature Dependent Voltage Exponent for Accelerated Life Tests (RCA)	242
D-2	Temperature Dependent Voltage Exponent for Accelerated Life Tests (Endicott & Zuellner)	243
D-3	Price and Reliability of Mica Capacitors	244
D-4	Price of Glass Capacitors	245
E-1	Failure Rate Compendium	248

FIGURES

Figure		
II-1	Composite Failure Distribution	28
II-2	Computer Cost and MTBF Versus Size ...	38
II-3	Computer Cost and Reliability	39
II-4	Integrated Circuit Failure Rate Versus Amount of Testing	55
II-5	Test Process for Semiconductor Device Procurement	60
III-1	Redundant Circuits	73
III-2	Redundant Circuits with Majority Voter	80
III-3	Asymptotic Availability of Redundant Computers (exponential service)	82
III-4	Availability of Redundant Computers (no service)	83
III-5	Probability of Failure During Peak Demand Period	89
IV-1	Comparison of Systems ($\mu = 1.0$)	97
IV-2	Comparison of Systems ($\mu = 0.1$)	98
IV-3	Comparison of Systems ($\mu = 0.5$)	99
IV-4	Comparison of Systems ($\mu = .025$)	100
V-1	Etched Board and Monolithic Integrated Circuit Layout	110
VI-1	Flowchart of Field Corrective Maintenance Process	115
VI-2	Relative Cost and Diagnosis Time Versus Parts per Module	126
VII-1	Asymptotic Availability of Redundant Computers (exponential service)	139
VII-2	Availability of Redundant Computers (no service)	140
VII-3	Comparison of Systems ($\mu = 1.0$)	142

VII-4	Comparison of Systems ($\mu = 0.1$)	143
VII-5	Comparison of Systems ($\mu = 0.5$)	144
VII-6	Comparison of Systems ($\mu = .025$)	145
A-1	The Exponential Distribution Function, $\bar{F}_c(t) = e^{-N\lambda t}$	158
A-2	Redundant Computer with Majority Voting	160
A-3	Availability of Redundant Computer-- No Service and Exponential Failure .	163
A-4	Availability of Redundant Computer-- No Service and Exponential Failure .	164
A-5	Availability of Redundant Computer-- No Service and Exponential Failure .	165
A-6	Availability of Redundant Computer-- No Service and Exponential Failure .	166
A-7	Availability of Redundant Computer-- No Service and Exponential Failure .	167
A-8	Availability of Redundant Computer-- No Service and Exponential Failure .	168
A-9	Availability of Redundant Computer-- No Service and Exponential Failure .	169
A-10	Availability of Redundant Computer-- No Service and Exponential Failure .	170
A-11	Redundant Computer with Imperfect Voters	171
A-12	Asymptotic Availability of Non- Redundant Computer (exponential failure and service)	176
A-13	Availability of Redundant Computer (transient phase, exponential service)	182
A-14	Availability of Redundant Computer (transient phase, exponential service)	183
A-15	Availability of Redundant Computer (transient phase, exponential service)	184

A-16	Asymptotic Availability of Redundant Computer (exponential service)	186
A-17	Asymptotic Availability of Redundant Computer (exponential service)	187
A-18	Asymptotic Availability of Redundant Computer (exponential service)	188
A-19	Asymptotic Availability of Redundant Computer (exponential service)	189
A-20	Asymptotic Availability of Redundant Computer (exponential service)	191
A-21	Asymptotic Availability of Redundant Computer (exponential service)	192
A-22	Asymptotic Availability of Redundant Computer (exponential service)	193
A-23	Asymptotic Availability of Redundant Computer (exponential service)	194
A-24	Comparison of Transient Availability, Exponential and Weibull Failure Distributions (non-redundant computer, exponential service)	199
A-25	Asymptotic Availability of Non-Redundant Computer, Exponential and Weibull Failure Distributions (exponential service)	201
A-26	Asymptotic Availability of Duplex Computers (flexible service)	204
A-27	Asymptotic Availability of Duplex Computers (limited Service)	208
A-28	The Multi-Processor	210
A-29	Asymptotic Availability of Multi-Processor (exponential service)	213
A-30	Asymptotic Availability of Multi-Processor (exponential service)	214
A-31	Asymptotic Availability of Multi-Processor (exponential service)	215
A-32	Asymptotic Availability of Multi-Processor (exponential service)	216
A-33	Asymptotic Availability of Multi-Processor (exponential service)	217

A-34	Asymptotic Availability of Multi-Processor (exponential service)	218
A-35	Asymptotic Availability of Multi-Processor (exponential service)	219
A-36	Asymptotic Availability of Multi-Processor (exponential service)	220
A-37	Asymptotic Availability of Multi-Processor (exponential service)	221
A-38	Asymptotic Availability of Multi-Processor (exponential service)	222
A-39	Asymptotic Availability of Multi-Processor (exponential service)	223
A-40	Asymptotic Availability of Multi-Processor (exponential service)	224

SYMBOLS

$f(t)$	Probability density function of T_f , the time to first failure.
$F(t)$	Distribution function of T_f , $\Pr[T_f \leq t]$.
$\bar{F}(t)$	$1 - F(t)$, $\Pr[T_f > t]$.
$g(t)$	Probability density function of T_s , the time to service.
$G(t)$	Distribution function of T_s , $\Pr[T_s \leq t]$.
$\bar{G}(t)$	$1 - G(t)$, $\Pr[T_s > t]$.
$h(t)$	Probability density function of $T_f + T_s$, the convolution of $f(t)$ and $g(t)$.
M	The number of redundant subsystems which comprise an entire computer (redundant computer).
"	The number of units per single computer (multi-processor).
m	$2m+1$ = the order of redundancy, $m=0,1,2$, (redundant computer).
"	The number of single computers in the system (multi-processor).
N	The total number of parts (assumed to be identical) in the entire system.
$P(t)$	The availability; the probability that the system is available (on) at time t .
P_∞	The asymptotic value of $P(t)$, $\lim_{t \rightarrow \infty} P(t)$
$r(t)$	The failure rate, $\Pr[t \leq T_f \leq t+dt T_f \geq t]$.
T_f	The time to first failure.
T_s	The time to service.
α	A parameter of the Weibull distribution.
Δ	A small time increment.

η	$M\lambda/N$
λ	The constant failure rate of a part whose failure distribution is exponential.
μ	The constant service rate of a service process whose distribution is exponential.

Chapter I

INTRODUCTION

This study began with a question: "Are contemporary computers reliable enough to be used in future ballistic missile defense (BMD) systems?" Several reasons suggest they might not be. Two major ones are: 1) the tremendous peak processing load for urban defense requires a computing complex of awesome size,[†] and 2) in the case of hard-point defense (e.g., hardened missile sites) the problem of servicing the computer in remote areas, even if its size were not comparable to the urban behemoth, still raises a serious doubt as to whether the data processor would be available for use if ever it were called upon.

It was subsequently decided to expand the question to include the entire subject of digital computer reliability, and to emphasize those large, ground-based computers where service is possible. The size of the processor and the nature of the service enter as parameters of a larger model. By providing estimates for the values of these parameters (i.e., by guessing the size of the machine and how to fix it when it breaks) the authors prevent the special problem (BMD data processing) from suffering by the generalization of the problem.

This decision has not been without dividends to the authors. It relieved them of the chore of actually estimating the size and logistic properties of any proposed

[†]For instance, the Univac development for Nike-X announced in Aviation Week, November 30, 1964.

BMD system, although this has partially been done for the hard-point defense case [1]. It does mean that the reader who wants to apply the results to a specific system must be able to say, for instance, "I've got about a 50,000 transistor computer made out of such-and-such quality parts and I'll probably spend an hour repairing it when it's down." There is much more to the story, but essentially the user (not so much the manufacturer) must have some idea of what is required to solve his computing problem and how he intends to provide the necessary service.

The reader should also know what level of reliability he desires. Whether the probability of the system being available (see p. 3) should be .90 or .99990 is left to the decision-maker, and the availability, as it is presented here, is only a measure which is functionally related to the many parameters of the system with no subjective value placed on it.

Some mention should be made at the outset about the use of certain terms. The word "reliability" appears in the title and continually in the text. Presently, a strict definition of reliability, along with other words, will appear, but it is not the authors' intent to aid in the proliferation of precise terms with which the reliability field already abounds. To this end, the word reliability has the usual colloquial meaning: being just a measure of whether "it works or doesn't."

One bias in viewpoint remains from the BMD beginnings. Namely, most of the effort in probabilistic analysis has been to ascertain if the computing system is on when it is needed. This differs considerably from asking, say,

about the fraction of all time that the computer is on. In other words, the assumption is usually made that the mission time is very much less than the lifetime of the system, and that the probability of successfully completing the mission given that the system is on at the start is almost unity.

The reader whose particular application won't permit this assumption, may, as a rule, answer the question, "What is the probability that the system is on at time t_1 and remains on until t_2 ?" by first finding the probability of being on at t_1 , then multiplying by the probability that no failure occurs in the interval (t_1, t_2) . This can be done for the majority of the examples given here because the relevant stochastic processes are Markov processes,[†] and in all cases the pertinent probability density function or distribution function will be given.

1. DEFINITIONS

Only three definitions are required for the work which follows.

- o Reliability. The probability that a device will perform its purpose adequately for the period of time intended under the operating conditions encountered.[‡]

[†] A Markov process has the property that the future of the process is only dependent on its present state and not on the time history of the process up to the present.

[‡] Radio-Electronics-Television Manufacturers Association, 1955.

- o Availability. The probability that the system can operate within the tolerances at a given instant of time.[†]
- o Interval reliability. The probability that at a specified time, the system is operating and will continue to operate for a duration, say x. The continued operation during the interval is, of course, to be achieved without benefit of repair or replacement [3].

Many more definitions might be useful, particularly those related to the fraction of time the machine operates. Most of these may be found in [2].

2. THE SCOPE OF THE REPORT

Below is an annotated guide to the main topics of this report.

Theoretical Performance. This is, in most cases, a computation and evaluation of the availability of a system. Knowledge of service and part failure distributions is assumed (and thus the word "theoretical"). The resulting systems distributions are computed and from these, the availability, $P(t)$, is computed. In many cases, the asymptotic availability, $P_{\infty} = \lim_{t \rightarrow \infty} P(t)$, is used instead of $P(t)$. The background and derivation is given in Appendix A instead of in the main text, and certain selected results are used in Chapters III and IV. The central problem is: Given parts of specified reliability, how should a computer be constructed to attain a desired system reliability? Subsequent chapters discuss achieving

[†]Hosford [2]. This is more formally known as "point-wise availability," and sometimes the "readiness."

this reliability through redundant techniques and multiple computers.

Component Reliability. Chapter II treats the problem of part failure. The word "part" is preferred over "component" since a component in this day and age may actually consist of a large collection of parts which, from a service standpoint, are indivisible. The origin of failures and their modes and mechanisms are examined and a large number of life-test statistics are analyzed.

This analysis provides the basis for conclusions about present and predicted reliability levels of parts, and how the reliability of a part relates to its price. From these data, particularly in Sections II-6 to II-8, it should be possible to extract an estimate of the reliability of a particular part on a simple, albeit very large, digital computer.[†]

The numerous details involved in part reliability appear in Appendices B, C, and D. The , the relationship between the behavior of a part and the stress to which it is subjected is discussed for semiconductor devices, resistors, and capacitors, respectively.

Appendix E presents a complete tabular compendium of the life-test failure statistics which were compiled during this study.

[†]By "simple," we mean that no circuit tricks have been used to increase the reliability, e.g., redundancy. Namely, a simple computer is assumed to consist of a large series-chain of fallible parts.

Circuit Design. Since it takes more than reliable parts to make a reliable computer, Chap. III considers the problem of obtaining reliable circuits. Chapter III presents an example of design philosophy which, if followed carefully, should go a long way toward guaranteeing a good design.

This chapter also discusses logical design and reports briefly on progress in reducing logical design errors by computer-aided techniques.

Chapter III introduces, as a tool to increase reliability, the technique of part redundancy. Good parts and good design are both necessary but unfortunately not sufficient to insure a reliable computer, so it is frequently necessary to resort to additional methods such as redundancy. The analysis appears in Appendix A, but the results are in this chapter, with some indication of what these methods can produce as a function of present and predicted part reliability.

Chapter III discusses means of failure detection. After all practical steps have been taken to obtain a reliable system, failures still occur; the problem then is to discover the failures.

Finally, III shows that the notion of a "failure" needs considerable refinement.

Multiple Computers. Many systems rely on the use of multiple computers to obtain reliability. Chapter IV introduces the use of spare computers to accomplish this goal, and evaluates the availability of such systems under different forms of service and system configuration. Chapter IV also discusses the relatively new concept of

the multi-processor, and compares the performance of the nonredundant, redundant, duplex, and multi-processor. Finally, it considers the problem of additional programming and hardware for the case of the multi-processor.

Systems Considerations. Chapter V takes up some of the factors that are important for matching the data processing system to its environment. Heading the list is the problem of writing error-free programs. Suggestions are made as to how this might be accomplished and further analysis is made of how such programs can be checked (no trivial problem in the case of real-time control processes) and corrected when errors are found.

Then such topics as interconnection and packaging reliability, performance in the face of extreme environments such as shock and radiation, preventive maintenance, and quality control are briefly discussed.

Maintenance. The central problem discussed in Chapter VI is discovering what to repair when the computer malfunctions. Current techniques in fault diagnosis and isolation will be presented and some estimates will be made of how effective and costly they are likely to be. This is a very important subject, since availability will be shown to depend strongly on the mean repair time, which can be drastically reduced if enough effort is expended on the problems of self-diagnosis.

Chapter II

THE PROCESS OF FAILURE

1. INTRODUCTION

A part will be defined here in the usual sense as the smallest replaceable element in a system. In other literature, parts are referred to as piece-parts, component parts, or components. Typical parts are: resistors, transistors, and integrated circuits. Note that the smallest field-replaceable element is a maintenance module, usually consisting of several parts, as described in Chap. III.

This study considers only electronic parts for two reasons.

- o Electromechanical devices used in computing systems are assemblies containing many non-standard parts for which there is no orderly body of performance and failure data such as exists for electronic parts.
- o The notorious unreliability of such electromechanical devices (typewriters, magnetic tape transports, etc.) seems to preclude their use in any function directly essential to the primary system mission, if extremely high reliability is required.

Furthermore, of the totality of electronic parts categories in existing and proposed computers, only a rather small subset will be examined in detail.

The major reasons for inclusion of a parts category for consideration in this study are:

- o High probability that the parts will be used in computers of the type under consideration.

- o High probability that the number of parts per system and predicted individual part reliabilities will contribute significantly to system failure rates.
- o High probability that there are or will be no better parts for the required part functions, within reasonable cost and time limitations.

Selected categories are listed below, as are many notable exclusions, with appropriate comments.

Silicon planar transistors--In ground applications, temperature is the most significant stress factor in semiconductor failure. Germanium exhibits asymptotic failure behavior at junction temperatures of about 150°C, whereas the critical temperature for silicon devices is more like 350°C. The cost of silicon transistors is still significantly higher than germanium, but progressively increasing silicon demand will narrow the gap. Also, it is not reasonable to compare cost of MIL-quality silicon to entertainment or commercial quality germanium.

The planar process is selected as it will most likely be the only surviving silicon fabrication process. Failure data for mesa devices can be included with planar data for extrapolations, and, with some caution, data on other transistor types may be included.

Silicon planar diodes--Essentially the same reasoning as above.

Monolithic silicon integrated circuits--For digital applications, the monolithic circuit is better than the hybrid thin-film circuit because it entails fewer nondiffused intraconnections, less handling in general,

and lower cost. Reduction to practice of laboratory methods for thin-film active device fabrication may shift the balance in favor of the monolithic circuits, but this hasn't happened yet.

Resistors--There are five types of resistors in general use. These are: composition carbon, carbon film (molded deposited carbon), tin oxide, metal film, and wirewound. Computers use few, if any, wirewound resistors, so this category may be ignored.

Capacitors--Dipped mica and glass capacitors will be considered for the range $1 - 10^4$ picofarads. Tabular paper or plastic capacitors cover the range $10^3 - 10^7$ picofarads, and tantalum and aluminum electrolytic capacitors are used for the range $10^6 - 10^{11}$ picofarads. A relatively new type, the multilayer ceramic capacitor, covers the range $10 - 10^6$ picofarads; lack of reliability information and probable economic inefficiency preclude consideration of this type.

In the most fundamental philosophical sense, a part has failed if, upon the future application of some combination of normally expected stresses, one or more of the parameters of the part would vary in such a way that the functional assembly containing the part would become incapable of performing its function.

This "philosophic failure" is as academic as the question of the sound of an explosion in the uninhabited desert, and a more practical definition might be:

A part has failed when, under some combination of normally applied stresses, one or more parameters of the

part vary in such a way that the functional group containing the part becomes incapable of performing its function.

From the above definition, it can be seen that it is difficult to establish a definition of part failure which is independent of the nature of the functional group (e.g., "circuit") containing the part. Circuits can be, and have been, designed to tolerate quite large variations in part parameters.[†] Also, circuits can be, and have been, designed to continue operating even when some number of parts have suddenly assumed limit values ("open" or "short" circuit).

Modifying an arbitrary definition of failure can significantly affect the relative merits of various types of parts. Consider the following example, based on manufacturers' published data [1]. Suppose very large samples of metal film and composition carbon resistors are placed on high-temperature load life test for several thousand hours. If failure is defined as "resistance variation more than 2 per cent from nominal," then nearly all of the composition resistors will "fail," but almost none of the film resistors. If failure is defined as "resistance variation more than 30 per cent from nominal," then the ratio of film resistor "failures" to composition resistor "failures" will be infinite, as none of the composition units will "fail," while a few of the film units will open.

[†]For instance, -100 to 1600 per cent from nominal I_{CBO} in transistors, -50 to infinite percentage on electrolytic capacitors, and h_{FE} of transistors, and ± 30 per cent on resistors.

At the part level, the usual practice is to classify relatively gradual or continuous parameter variations with stress (time included) as "degradation" and reserve the term "failure" (usually with the superfluous adjective "catastrophic") for the relatively rapid or discontinuous passage of a parameter to a limit value (open, short) or a value outside some statistically predicted bounds for degradation behavior.

To avoid the qualitative judgment implied by the words "relatively rapid or discontinuous," they may be deleted from the definition of failure. This change does, however, introduce "inadequate prediction of degradation bounds" as a failure mechanism.

Whether in operational use, under test, or on the shelf, a part is characterized by a set of parameters p_i , $i=1,2,3,\dots,m$, and subjected to stresses s_j , $j=1,2,3,\dots,n$, where time is included as a stress. The parameters are functionally related to the stresses by writing $p_i=f_i(s_j)$, where the form of the f_i may be quite complicated for even the simplest parts and in many instances unknown. In the degradation range, which might be loosely defined by $H_i \geq p_i \geq L_i$, where the H_i and L_i may be constants or functions of the s_j , efforts are usually made to find approximate forms $p_i \cong \sum f_{ij}(s_j) \cong \sum a_{ij}s_j$ where a_{ij} are coefficients of the first-order terms in the series expansion of the f_{ij} , and represent such often-used coefficients as "temperature coefficient of resistance," and "voltage coefficient of resistance." When the f_{ij} are not so tractable, bounds may be given, such as variation at 100 per cent rated load, 1000 hr, +2 per cent to -4 per cent maximum. For practical

application, bounds of this type are essentially statistical limits for any part which may be used.

The point to be made is that good information relative to degradation of all parameters is available for most types of parts which are under consideration here. If the predicted degradation behavior of a certain type of part is inadequate for the proposed application after all circuit efficiency tradeoffs have been considered (Chap. II), a better part must be selected.

More detailed examination of the meaning and origin of "failure" is necessary because the distinction between degradation and (catastrophic) failure is quite arbitrary. Setting $H_i = +\infty$, $L_i = -\infty$ categorically reduces the failure rate to zero for all components, but leaves the circuit designer very little to work with. Consider, for instance, the following cases of failure.

Failure occurs if a combination of s_j exists such that one or more $p_i > H_i$ or one or more $p_i < L_i$. This situation in turn has several interpretations. If the combination of s_j was expected (within absolute bounds, or statistically predicted) then either the limits H_i , L_i were assigned with the knowledge that the failure possibility existed, or in ignorance. Assignment of ± 2 per cent limits to a film resistor with the knowledge that an "open" can occur is nevertheless reasonable, as it would not be practical to assign -2 per cent, $+\infty$. If the combination of s_j was unexpected, then either the absolute bounds were wrong or a statistically unlikely combination occurred.

Failure also may occur if some previously unknown or unexpected stress appears (some s_k , $k > n$) or some new functional relationship appears between the expected s_j and the p_i .

Finally, there may be some truly random failure mechanisms, even some which do not obey known causal relationships.

Currently the problem of reliability of electronic equipment is approached either through statistical analysis and synthesis or physical reasoning.

The statistical approach is essentially a carry-over into electronics from other fields, notably mechanical engineering. Basic to this method is the assumption that the gathering of sufficient field and test data on samples of parts manufactured by some relatively constant process permits fitting some mathematical functions to failure behavior. The classic example yields the three-phase failure rate curve (sometimes called the "bathtub," probably by the same individuals who refer to the normal density function as the "bell") which is a composite of a decreasing failure rate for time near zero (early failure or "infant mortality"), a constant failure rate for some intermediate period and finally, an increasing failure rate at some later time (the process of "wearout").

Wearout failures result from physical or chemical changes with time, temperature, and other stresses, which either cause a parameter to exceed a degradation bound (as in oxidation of a metal film resistor) or to go suddenly to a limit value (as in the work-hardening and eventual breakage of a relay armature spring).

The origin of failures whose failure rate is constant is not so clear. In mechanical systems, random failures have been said to originate from simultaneous combinations of stress, randomly and independently occurring, which exceed the strength of the part. This implies that either all stress limits were not known, or the part was intentionally not designed to stand all combinations of stress. In electronic work, it is usually possible to design parts which simultaneously tolerate all anticipated maximum stresses. If this is done on an absolute basis (so-called "worst-case" design), then only ignorance of the limits would permit random failure from stress. If design is done on a statistical basis, a random failure mechanism exists.

There are true random mechanisms affecting very small electronic parts, notably semiconductor devices, at the atomic and subatomic level. Various types of crystal defects may be caused by internal statistical behavior as a function of time and temperature, and the ambient radiation environment. Under normal circumstances (e.g., 55°C, ambient temperature and no recent nearby nuclear explosion) mathematical estimates of these effects put them several orders of magnitude below various macroscopic effects as failure modes [2].

Early in the history of reliability analysis, when the state of the part manufacturing art was not so well advanced, straightforward life-testing of samples of parts produced sufficient failures to permit reasonably good curve-fitting and other estimates of failure behavior. In six months or so, a sample of several hundred incandescent lamps could be run until nearly all had failed.

By contrast, a sample of a thousand semiconductor devices might be run for a year with two, one, or even zero failures.

Two approaches to the problem of insufficient data have been used--larger sample sizes and accelerated life testing. Ever-increasing sample sizes produce more failure data, but at proportionate cost increase. The hypothesis underlying accelerated life testing is as follows: if the incidence of failure and the level of a particular stress are functionally related in some well-behaved manner, the failure rates observed at two or more high stress levels may be used to predict failure rates at lower stress levels by extrapolation. Many researchers have shown excellent correlation between observed failure data and simple functions of absolute temperature for resistor and semiconductor devices [3]. Others, however, claim that the observed relations hold only at the higher temperatures, and that the extrapolation back to lower temperatures is meaningless, as the high-temperature failure modes (e.g., oxidation and phase changes) are virtually non-existent, while other modes exist which are not acceleratable by increased ambient temperatures.[†]

With respect to the procedure of increasing sample sizes to obtain useful numbers of failures in reasonable times, the occurrence of certain new failure modes makes

[†] A major argument against step-stress (accelerated) testing has been the discrepancy between results obtained with power-stressing and high temperature (unpowered) stressing [4,5]. Recent information [6] indicates that incorrect determination of junction temperature in power-stress tests may account for the observed differences.

it doubtful whether extrapolation to failure characteristics of single units is valid. In other words, there is reason to believe that the failures do not constitute a Markov process and that, as we test for a longer time, new failure modes occur. It would certainly be questionable, therefore, to draw conclusions about the behavior of a single device over 100,000 hr, based on observations of 100,000 devices for one hour.

Another consequence of the "brute-force" testing method is that the most and best information is available on the oldest, i.e., most obsolescent, parts.

All of the above considerations have led to the evolution of the physics-of-failure approach to reliability prediction. This approach lists and classifies all significant part failure mechanisms, and establishes, on a theoretical basis, the functional relation to stress. Determination of possibly significant mechanisms originates from observed failure modes or from pure physical reasoning. Where possible, functional relations are correlated to observed data to confirm hypotheses. If individually structure-related and material-related failure mechanisms are isolated, and functional stress relations established, the same relations may be carried over to different parts produced by similar processes, without repeating extensive tests.

Cautious application of physics-of-failure techniques can actually identify mechanisms which are amenable to accelerated testing, which in turn permits verification of hypotheses, at least for the well-behaved cases [7].

2. FAILURE MODES AND MECHANISMS

Five terms relating to failure behavior need definition.

Indication--The external observation that a part parameter has changed to a value outside the established degradation bounds, often to a limit value (e.g., open, short, no output).

Mode--The internal occurrence which causes the indication.

Mechanism--The physico-chemical process underlying the mode.

Stress--Any characteristic of the environment which causes or allows the mechanism to proceed.

Origin--That aspect of the materials and processes of fabrication of the part which allows a combination of stress and mechanism to result in a mode.

Application of the terms is exemplified in the two cases of integrated circuit failure shown in Table II-1. These two examples, taken by themselves, show that the origins of failure fall into two classes.

- o Fundamental design inadequacies resulting from characteristics and limitations of materials, processes, procedures, equipment, and personnel.
- o The degree to which potential capabilities of the design are realized in actual fabrication (production engineering, process control, quality control).

A third factor which affects failure in a "post facto" sense is removal of actual or potential defectives,

Table II-1

EXAMPLES OF INTEGRATED CIRCUIT FAILURE

	Case 1	Case 2
Indication	Open input circuit	High reverse current
Mode	Thermocompression bond off at aluminum pad	Excessive moisture in can
Mechanism	Formation of inter-metallic compound ("purple plague")	Ionic conduction on surface of semi-conductor
Stress	Temperature, time	Ambient air pressure, humidity, time
Origin	Use of bimetallic gold-aluminum system in presence of silicon	Faulty weld, permitting air leakage into can

regardless of origin, before use in equipment or before the equipment is declared operational. The processes involved include:

Inspection and non-destructive test--Visual, electrical, radiographic, infra-red, mechanical, package leakage, and other tests, conducted either on a sampling or 100 per cent basis.

Destructive test or test-to-failure--Usually step-stress tests using temperature, power, vibration, or acceleration. Of necessity these tests are conducted on a lot sampling basis only.

Screening and burn-in--Subjection of parts to one or more stresses carefully selected so as to greatly accelerate decreasing failure rate mechanisms without significantly affecting any constant or increasing rate mechanisms. Examples are high temperature bake, temperature cycling, acceleration, and operation in test circuits. Burn-in is usually conducted on a 100 per cent basis.[†]

Debugging--Replacing early failures as they occur during equipment checkout--a necessary procedure which operates on a 100 per cent basis and is essentially a form of burn-in. Debugging is inferior to laboratory burn-in, however, as the parts in actual equipment, hopefully, are derated (understressed) while well constructed burn-in tests often apply carefully selected overstresses. Physically small (e.g., airborne) computers may be burned-in by operating in high stress, usually high-temperature, environment, but this is still an undesirable procedure, as the least-derated parts would have to be grossly overstressed if the most-derated parts are to be significantly stressed. Furthermore, practical aspects of checking out a large-scale, ground-based computer at, say, 125°C for a week, leave something to be desired.

[†]"Burn-in" is the process of actually operating the part in a controlled environment for a length of time that, hopefully, will weed out parts that were destined for early failure. Effectiveness of the burn-in procedure has been demonstrated by the Apollo project in which burn-in procedures requiring approximately 15 days have produced 0.2-0.3 per cent early failures, whereas 17,000 parts surviving burn-in have operated for an average of 120 days per part (49 million unit hours) with no subsequent failures.

To date, most large-scale reliability testing and theorizing has been conducted in support of airborne, spaceborne, missile, and shipboard systems. These are usually small computers designed for high probability of survival over mission times, or times between maintenance, of minutes to a few hundred days. The concept of a "useful life" or "design life," of several years, in the few cases where applicable, is nullified by the probability of technological obsolescence. Most manufacturers of large commercial computers do not support massive electronic parts reliability programs, for several reasons.

- o Reasonable extrapolation can be made from military data.
- o Parts manufacturers supply "computer grade" parts made by processes similar or identical to those used in military parts, but without the qualifying paper burden.
- o In-house parts reliability programs cost too much.
- o Electro-mechanical peripheral equipment failure overshadows electronic parts failure in most commercial systems.
- o Design life is most likely limited by technological obsolescence.

Exceptions to the above might be noted for real-time industrial control systems of hazardous processes (steel mills, refineries), and communications processing where errors and delays incur significant costs (stock market quotations systems, on-line teletype data processors, subscription television billing computers). Other large-scale ground-based systems requiring long design life are the military warning, strategic, and tactical machines.

The need for long life in military computers comes about primarily because the military purchases the computer and the manufacturer makes no provision to offer the military a "new model" every three or four years.

The outstanding civilian exception in the realm of non-military reliability is, of course, the Bell Telephone System. Typical of the fundamental telephone company attitude toward reliability is the design life target for solderless, wire-wrapped connections used in central exchanges: 40 years. This significant difference in required lifetime, plus the dynamic state of parts technology, calls for critical review of failure theory and data [8].

At the parts level, the most significant question is-- are there any increasing failure rate mechanisms in any of the part-types that may be selected; and if such mechanisms exist, do they exhibit asymptotic or sharply peaked behavior at times within the design life of the computer? Consider the (unlikely) hypothesis that operation of the system depends on the operation of a large number of incandescent indicator lamps. The immediate reaction of a design engineer might be to suggest use of the available (at a price) 10,000-hr lamps, or the (at a higher price) 50,000-hr lamps. But, from the economic standpoint of a lamp manufacturer, lamp life should be normally distributed about the advertised value, with a very small variance. This, coupled with the fact that a reasonable 10-year design life contains 87,600 hr, indicates that the intuitive response to design requirements of the subject system may be grossly in error.

If empirical evidence or prediction based on theory shows existence of increasing failure rate mechanisms, there are still two possible mitigating situations which allow us to live compatibly with such mechanisms: the moments of the failure distributions are such that reasonable systems readiness may still be achieved, or the mechanisms are distributed non-uniformly among the population of parts. The first situation is straightforward; if a "wearout" mechanism makes itself known at 50 years, as might be the case for solid state diffusion (see Appendix B), it will probably not significantly affect system readiness over the first ten years. As an example of these statements, Table II-2 shows the mean failure rate over the first ten years of operation for parts whose failure distribution is normal with given mean and standard deviation.[†]

Evidence exists that certain increasing failure rate mechanisms have origins which are not uniformly distributed among the part populations. Certainly failures like case 2, in Table II-1, resulting from obvious manufacturing defects, are in this class.

The effect of certain contaminants (metallic particles, water vapor, etchant residues) in transistor, diode, and integrated circuit packages is worthy of consideration in this respect. It may be that contaminants are only included in a small percentage of the population. If so, and if adequate quality control maintains or improves the

[†]Actually, a truncated normal distribution, but the left tail, for $t < 0$, is negligible in all but a few cases.

Table II-2
NORMALLY DISTRIBUTED WEAROUT

Mean Wearout Life, (years)	Standard Deviation (years)	Mean Failure Rate over 10 years (%/10 ³ hr)
20	5	.026
20	10	.181
30	5	.00037
30	10	.026
40	5	<10 ⁻¹⁰
40	10	.0015
50	5	<10 ⁻¹⁸
50	10	.000037

status quo, then the contribution of wearout by contaminant effects to overall failure rate applies in the same small percentage. Contaminants may, however, unavoidably be distributed throughout the entire population.

A uniform distribution of contaminants may cause a severe wearout hazard. If the distribution is nonuniform (normal, say), the hazard may be less severe, particularly so if the cause-effect relation between the contaminant concentration and the failure mechanism is nonlinear or discontinuous (exhibits a threshold).

The preceding discussion was hypothetical. We will now attempt to pin down actual failure behavior, although, unfortunately, the available evidence is meager, and the many conflicting interpretations only confuse the situation.

3. MODELS OF FAILURE BEHAVIOR

Consider first the ideal device manufactured exactly in accordance with its design. For a single "ideal" part (metal film resistor, say), if there is one, or at most a few, dominant failure mechanisms, life testing plus the physics of failure (writing of a theoretical/empirical equation for each mechanism plus isolated parameter determination of each mechanism) can yield useful information.

For a complex part (e.g., transistors, or integrated circuits) many mechanisms may act simultaneously and their degree of contribution to the probability of failure may vary strongly with instantaneous stress values. Further, the nonhomogeneous nature of the materials and possibly less-well-understood behavior makes the relation between physics-of-failure equations and reality even more tenuous.

Consider now the real, as-built device, in which all the statistical variability of the materials and manufacturing processes has been superimposed on the parameters of the basic design.

It can be categorically stated that, in early life, for complex parts (and even a film resistor or mica capacitor may be complex in this sense), some form of decreasing failure rate will be observed which is in no way correlated to predicted behavior of the ideal design. The rate of decrease may be correlated to the vendor's

name [9,10],[†] relative newness of the product or process, or even the material of an assembler's skirt and motion of that portion of the assembler inside the skirt relative to her chair [12]. This behavior is essentially the failure of "genetic defectives" under normal stress or (in the case of burn-in) intentional, controlled overstress.

It must be true that, in later (be it decades or centuries) life, all parts will show increasing failure rate ("aging" or "wearout") mechanisms. At any non-zero temperatures, solid state diffusion affects all devices, semiconductors most significantly. With applied voltages and currents, drift and migration occur. With any reactive elements in the environment, oxidation, electrolysis, and similar actions take place. And with any cycling, stresses, work-hardening, fatigue, crystallization, and defect propagation occur.

In the period when the early failure rate has decreased to zero and the wearout rate is still insignificant, there may be some very low-level, truly random (non-causal) failure mechanisms operating.

Over the duration of all life tests performed so far, and in the opinion of nearly all authorities (e.g., Refs. 13, 14, 15), the decreasing failure rate behavior dominates. Certainly this is true for tens of thousands of hours, and probably through ten years.

The few dissenting commentaries point out existence of shorter-term wearout mechanisms (intermetallic phase formation, anomalous compound formation, inversion due to seal leakage, oxide regrowth at windows) but the effects

[†]Also, see p. 170 of Ref. 11.

of all of these are apparently present in only a small percentage of the total population. In fact, the very presence of these mechanisms creates the long "tail" of the decreasing failure rate curve. To reiterate: the total part population shows a decreasing failure rate because, and only because, various controllably small subgroups show initially increasing failure rates of various forms until every member of the subgroup has failed, at which time the failure rate of the entire remaining population effectively decreases.

Adequate shakedown (checkout) periods, or much preferably, carefully designed burn-in procedures can locate the operational time origin far down on the decreasing rate slope without intruding on the upslope of the long-term true wearout mechanisms.

After burn-in, the constituents of the composite failure-rate curve are mostly post-modal tails of the subgroup functions, plus any very-low percentage, high variance, high-mode subgroup functions, plus the background true random rate, plus the premodal tails of the true wearout functions.

Figure II-1 shows a synthesis of a failure distribution function which was constructed under the following hypotheses.

- o Wearout mechanisms affecting the entire part population have a mean between 10 and 100 years.
- o There are six wearout mechanisms operating on subsets of the population, with means of 30 to 30,000 hr.
- o The mechanism with the highest mean (about 30,000 hr) affects the smallest (about 1.7 per cent) percentage of the population and mechanisms with progressively lower means affect higher percentages, with the 30-hr mechanism affecting 10 per cent of the group.

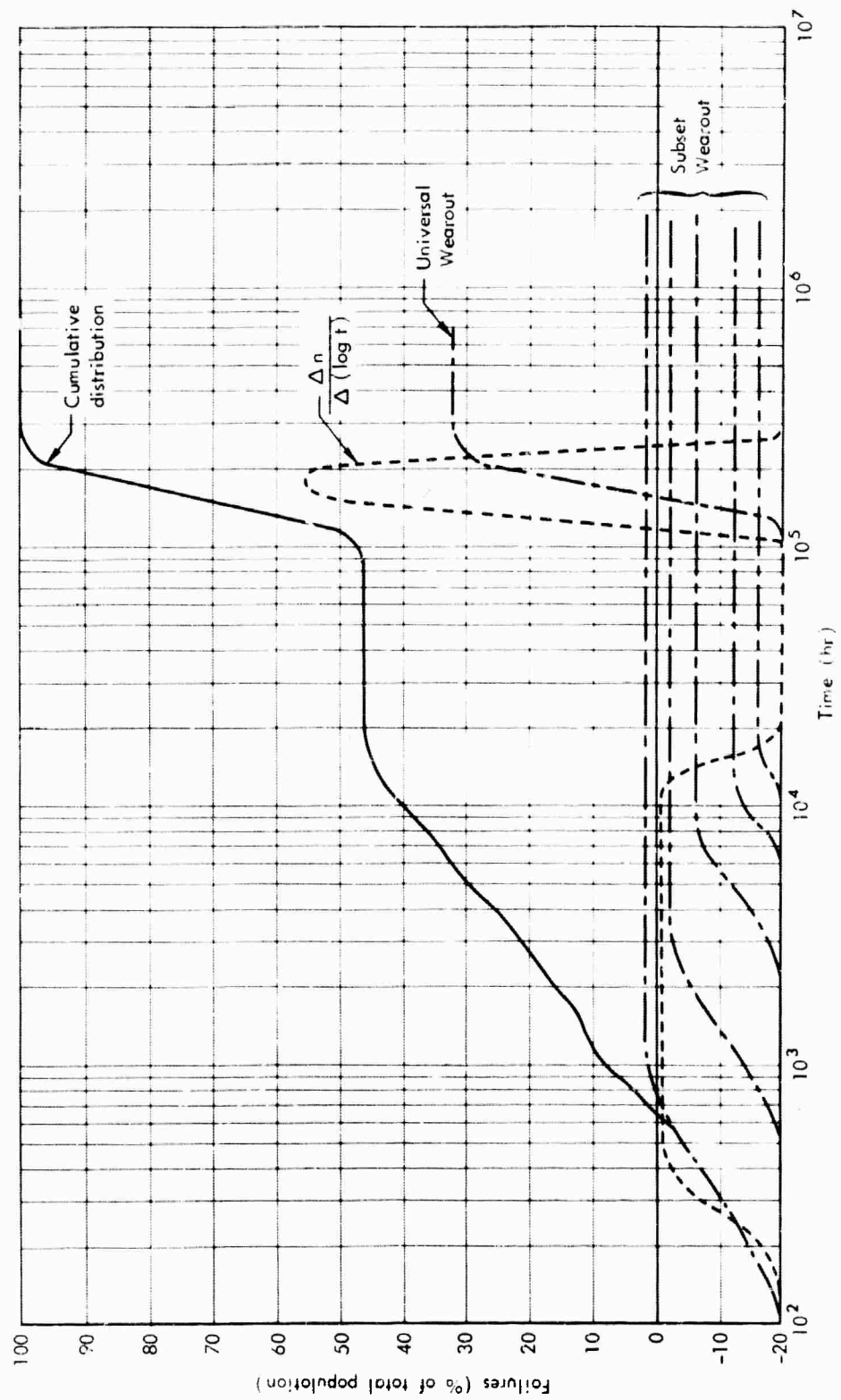


Fig.II-1—Composite failure distribution

The composite curve shows a very short increasing failure rate period, followed by a decreasing failure rate to about 2×10^4 hr, at which time the rate has decreased to zero. There is then a "golden age" from 10^4 to 10^5 hr (about ten years) where the failure rate remains zero, or more properly is reduced to the true random background rate. At 10^5 hr, the universal wearout mechanisms begin to take effect.

The approximate failure rate vs. time for Fig. II-1 is tabulated in Table II-3.

Table II-3

FAILURE RATE VS. TIME FOR SYNTHETIC FAILURE DISTRIBUTION
OF FIG. II-1

Time (hr)	Failure Rate (%/10 ³ hr) (over random background)
15	170
30	250
70	105
150	40
300	17
700	5.8
1,500	4.5
3,000	1.5
7,000	0.3
15,000	0.05
30,000	0.00
70,000	0.00
150,000	0.17
300,000	0.13
700,000	0.04

The time scale is to be considered as "equivalent unstressed hours," where legitimate accelerating processes exist. The effect of burn-in may be noted by "throwing away" that portion of the cumulative distribution to the left of the equivalent time represented by the burn-in process.

Faced with evidence of multiple early-wearout mechanisms in semiconductor devices, plus the likelihood that they are much less-well-behaved in reality than in the example of Fig. II-1, it is difficult to see how efforts to fit various popular distributions (notably the Weibull and the log-normal) can succeed. Some attempts at fitting are given in [16] and [17], and the abandonment of an heroic attempt is stated in [15]. Descriptions of the distributions used in the fitting game may be found in [18].

The question remains--how can system reliability be computed without theoretical or empirical failure distribution? There seems to be but one rational answer: When there is no additional information, choose a distribution with a constant failure rate (this is uniquely the exponential distribution[†]). This constant failure rate, λ , may be obtained by computing a weighted average of the true failure rate over the design life and defining λ to be this number, or defining λ to be piecewise constant over the interval of interest. These constant pieces should be estimated by the educated combination of all the available processes--physics of failure, brute force life tests, accelerated step-stress tests, and engineering intuition.

[†]If T_f is a random variable denoting the time to failure, then T_f is exponentially distributed if $\Pr[T_f \geq t] = e^{-\lambda t}$, $\lambda = \text{constant}$. See Appendix A.

Eliminating unreasonable figures from the mass of available data purporting to indicate inherent failure rates requires engineering intuition (plus, possibly, some detective work). There are two types of offenders--unreasonably high failure rates and suspiciously low failure rates. Causes of unusually high failure rates include:

- o Use of field data which includes secondary failures, "homicides," i.e., destruction of parts resulting from human error, and replacement of parts which were not truly defective. Homicides are particularly significant. On one large computer (450,000 semiconductor devices) the transistor failure rate was $.0013\%/10^3$ hr in one quarter and $.114\%/10^3$ hr in the preceding quarter. Inquiry produced the following approximate statement "Oh yes, accidents happen. Just last month a man dropped a probe and took out 100 modules." (In this example, then, human error introduced about two orders of magnitude increase in reported failure rate.)
- o Use of data on a new product for which the quality control process has not had time to operate.
- o Use of data from tests conducted at stresses significantly higher than those anticipated in the application.

Causes of unusually low failure rates include:

- o Extrapolation from accelerated or step-stress tests where it is not clear that this is physically legitimate, or such extrapolation using some empirical formula in a range far from that in which its parameters were determined.
- o Discarding certain failures from the statistics as errors in fabrication or for other reasons, unless it is clearly proven that the origin of the defects has been removed.

4. ESTIMATES OF NON-CONSTANT FAILURE RATES

The best available evidence and opinion indicates that, for the parts under consideration, there are no increasing failure rate mechanisms which will have significant effects in the first few decades of part life.

A considerable body of opinion and evidence indicates decreasing failure rate behavior of many part-types, notably semiconductor devices and certain capacitors. It is not clear that the decreasing failure rate behavior (observations) is indicative of the existence of corresponding decreasing failure rate mechanisms. As stated earlier, superposition of various subgroup short-term mechanisms probably produces the observed group behavior.

Nevertheless, if the group behavior is as observed, and if screening, burn-in, and quality control cannot remove "genetically defective" individuals, a decreasing failure rate may have to be seriously considered in predicting system reliability, unless the rates settle to some constant or near-constant background value in a small fraction of the design life (e.g., one year for a ten-year system).

Reference 15 gives one of the largest collections of data gathered for a single investigation of non-constant failure rate behavior. Eight manufacturers contributed life test data on 10,300 transistors of 24 types. Maximum test time was 1000 hr for all but two types. Most tests were at reasonably high temperature or dissipation levels, and definitions of failure were arbitrary limits of two ranges--initial and end-of-life.

Failures were logged at various elapsed-time intervals. The Weibull failure distribution was assumed, and attempts were made to compute the parameters α and λ of the Weibull distribution and the associated confidence limits.[†] After discarding high-power and unijunction transistor data, and those cases for which insufficient failures occurred, Table II-4 gives what meager information remains.

Table II-4
ESTIMATED PARAMETERS FOR WEIBULL FAILURE DISTRIBUTION

Type	Junction Temperature °C (Operating or storage)	Weibull Parameters ^b			
		Initial Limit Failures		Life Test Limit Failures	
		α	$1/\lambda$	α	$1/\lambda$
2N652A	100 sto	1.15	10.0	1.00	60.3
2N652A	100 op ^a	0.66	16.4	-	-
2N705	300 sto	0.15	45.0	0.23	77.0
2N705	100 op ^a	0.53	55.0	-	-
2N718A	200 op ^a	-	-	0.29	61.0
2N744	175 op ^a	0.56	6.62	-	-
2N962, 964	100 op ^a	1.00	50.4	0.55	66.7

^aEstimated from dissipation and thermal resistance.

^bWith t in thousands of hours.

[†]If T_f is the time to failure, then T_f is Weibull distributed if $\Pr\{T_f \geq t\} = e^{-\lambda t^\alpha}$. The failure rate for the Weibull distribution is $r(t) = \lambda \alpha t^{\alpha-1}$.

For the data as given, enormously large failure rates result. Some effort is necessary to extrapolate observed values to more reasonable operating conditions. Assuming 55°C ambient and 15°C rise gives an operating junction temperature of 70°C for the proposed system. Acknowledging the hazards of the process, an attempt was made to produce an Arrhenius extrapolation from the data, assuming

$$\log \lambda(T_2) = \log \lambda(T_1) - k \left(\frac{1}{T_2} - \frac{1}{T_1} \right)$$

Values of k obtained from Refs. 2, 5, 17, 19, and 20 were as follows: 1.47, 2.20, 4.44, 4.99, 5.21, 5.70, 6.20, 6.23, 19.1, for λ in %/1000 hr and reciprocal temperature in $1000/T^\circ\text{C}$.

The results of this extrapolation to 70°C resulted in the following "best 1965" estimates for the part failure rate, if a Weibull distribution is assumed.[†]

$$\text{Transistor: } r(t) = .005t^{-0.4}$$

$$\text{Diode: } r(t) = .0025t^{-0.4}$$

$$\text{Resistor and Capacitor: } r(t) = .00017t^{-0.4}$$

[†]The validity of these failure rates is marginal at best. Both the original data and the extrapolations are suspect; but since the subject of estimating decreasing failure rate keeps arising, the authors decided to include these data.

5. STORAGE LIFE VERSUS OPERATING LIFE

If there is only a periodic demand for a system, it is essential to consider the effect on reliability of putting the system in some standby condition between operating periods. "Standby" might mean a completely de-energized state or a carefully designed condition where the parts are subjected to some optimal environment. Effects of environment, in this context, on various part-types will first be considered.

Composition resistors [1]--Humidity and voltage degradation effects are largely reversible. Operation at at least 1/10 rated dissipation, or in a controlled-humidity environment, will minimize humidity effects. Temperature effects are partially reversible, while load-life effects are relatively permanent. The optimum standby condition would be: low ambient temperature and humidity, and zero power dissipation.

Film resistors--Degradation mechanisms are enhanced by power dissipation and temperature. However, a large number of power-temperature cycles might increase the probability of "open" failure. Optimum standby condition: probably left on, if well-derated.

Capacitors (all non-electrolytic)--Life is a sensitive function of voltage and temperature. Optimum standby condition: zero volts and low ambient temperature.

Electrolytic capacitors--Voltage, surge current, and temperature decrease operating life, but some forward

voltage is required to prevent de-forming. Optimum standby condition: low ambient temperature, forming voltage applied through current-limiting resistor.

Semiconductor devices--Early and long-term mechanisms are enhanced by voltage, current, power, temperature, and humidity. Cycling also may have significant effects, and sensitivity to over-voltage transients is extreme. Optimum standby condition: low ambient temperature and humidity, zero power dissipation.

Although the above discussion clearly defines an optimal standby condition for a system, there are some very strong arguments against standby operation. If the system must be energized daily or more often for self-check purposes, the potential hazard of stress-cycling and uncontrolled transient damage must be very carefully evaluated.

There is a considerable economic justification (and pressure) to operate the system for routine computation when it is not performing its primary role. Also, most systems (e.g., ballistic missile defense) are in a surveillance mode and must therefore be on, although operating well below capacity.

If the temperature, humidity, and power derating of the parts are well-controlled, there is little absolute difference in system failure rate of the machine when operating and when it is on standby--perhaps a factor of two, at most. If this is weighed against the value of the system capability for peripheral tasks, and the risk of cycling or transient effects, there does not seem to be any clear-cut advantage in the standby mode.

6. COST-RELIABILITY RELATIONSHIPS

Functional relationships between inherent reliability and cost are usually difficult to obtain for the following reasons:

- o Demonstrated reliability, for parts already on test, steadily improves with the passage of time, unless and until
- o A single failure occurs, which instantaneously and dramatically increases the proven failure rate, certainly with no accompanying change in product cost.

It is nevertheless reasonable to assume a correlation among reliability, quality, and cost. Furthermore, attempts by manufacturers to qualify parts to various reliability levels must involve two major aspects of quality improvement:

- o Use of basic research, along with destructive test and field failure data, to modify and improve materials and processes;
- o Maintenance of quality of materials and uniformity of processes.

The cost/reliability/complexity graphs shown in Figs. II-2 and II-3 (which will be described at the end of Sec. II-6) were prepared for three approaches to system parts procurement:

- o Buy good commercial-industrial (computer-grade) parts and use as-received;
- o Buy as above but perform limited in-house screening and burn-in;†

†Such in-house tests include: Transistor-measure and record h_{FE} and I_{CBO} , then bake, temperature cycle and centri-

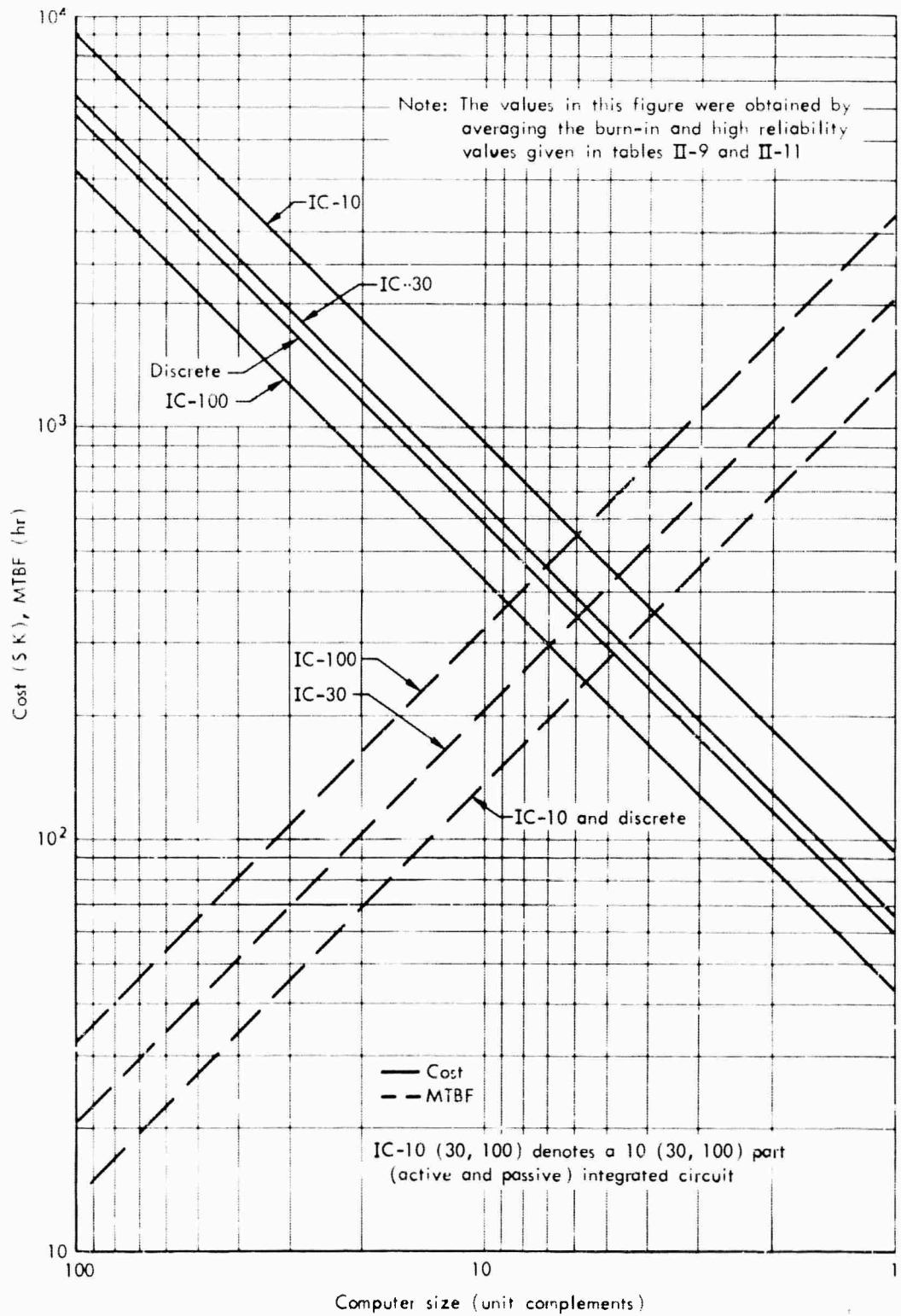


Fig. II-2—Computer cost and MTBF versus size

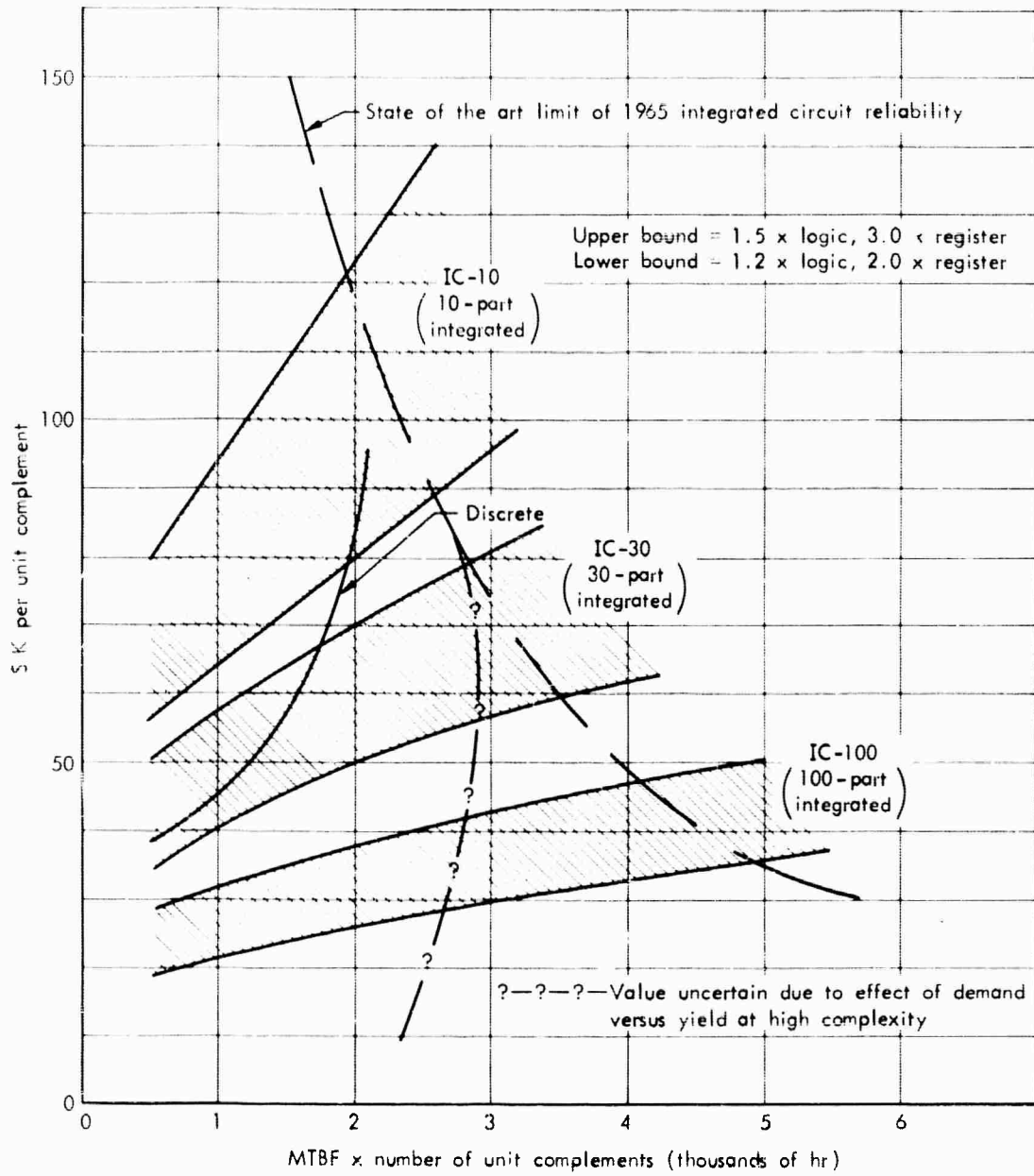


Fig.II-3—Computer cost and reliability

- o Buy high-reliability parts and use as-received. "High-reliability" here assumes some procedure such as two stages of 100 per cent screen and burn-in, sampling tests for mechanical and environmental stresses, long-term life tests, and a quality assurance program.

Individual part costs, in reasonably large quantities, are estimated below in Table II-5, with estimated failure rates in %/1000 hr in parentheses.

Relations between the failure rates permit expressing the system complexity in "equivalent transistors" and use of only the transistor failure rate in reliability calculations. The complexity is defined as

$$\text{Complexity} = T + D/2 + (R+C)/30 ,$$

where T = number of transistors, D = number of diodes, R = number of resistors, C = number of capacitors.

Approximate complexities of some existing systems are shown in Table II-6 for orientation.

The percentage of each type of part in an average computer is estimated in Table II-7.

If a unit parts complement is taken as 10,000 transistors, 15,000 diodes, 15,000 composition resistors, 5000 film resistors, and 5000 mica capacitors, the complexity

fuge, measure and record h_{FE} and I_{CBO} again and reject any failed or deviant units; Diode--same as above, measuring V_{fwd} and I_{rev} ; Integrated circuit--as above, measuring selected transfer relations; Resistor--measure R, temperature cycle and measure R; Capacitor--measure C and I_{leak} , and apply simultaneous voltage and temperature stress, measure C and I_{leak} .

Table II-5
COSTS FOR THREE GRADES OF PARTS

Item	Computer Grade (\$)	Computer Grade with User Screen and Burn-in (\$)	High- Reliability (\$)
Transistor	3.00 (.01)	3.35 (.005)	4.50 (.003)
Diode	0.50 (.005)	0.65 (.0025)	1.00 (.0015)
Comp. Resistor	0.04 (.0003)	0.08 (.0002)	-
Film Resistor	0.10 (.0003)	0.15 (.0002)	.50 (.0001)
Mica Capacitor	0.04 (.0003)	0.08 (.0002)	.35 (.0001)
Int Ckt, ^a			
10 parts ^b	10.00 (.02)	12.00 (.009)	15.00 (.005)
30 parts	20.00 (.04)	24.00 (.02)	30.00 (.009)
100 parts	40.00 (.07)	48.00 (.03)	60.00 (.015)

^aIntegrated circuit--price estimated for mid-1965; availability and price of more complex circuits somewhat uncertain.

^bThe parts are both active and passive.

Table II-6

COMPLEXITY OF EXISTING SYSTEMS

System	Complexity in Equivalent Transistors
FSQ-32	383×10^3
FSQ-31V	274×10^3
CDC-3600	97×10^3
CDC-1604A	82×10^3
Univac 1107	78×10^3
Burroughs B-5000	67×10^3
Honeywell H-1800	49×10^3
Honeywell D-825	41×10^3
SDS 9300	35×10^3
USQ-20	30×10^3
IBM 7090/44	26×10^3
GE 215/225/235	22×10^3

Table II-7

PART PERCENTAGES OF AN AVERAGE COMPUTER

Part	Percentage
Transistor	20
Diode	30
Composition resistor	30
Film resistor	10
Mica capacitor	10

and cost for each of the three procurement policies listed on pp. 38-43 are given in Table II-8.

The failure rate, appearing in Table II-9, is the product of the number of equivalent transistors by the transistor failure rate in each procurement category.

Evaluation of configurations using integrated circuits of various complexities requires some cautious interpretation. The results below are based on the following assumptions.

- o One can integrate 80 per cent of the system. The remaining 20 per cent, such as line and memory driver circuits, remains discrete.
- o The integrated portion is $\frac{2}{3}$ logic-type circuits and $\frac{1}{3}$ flip-flop or register-type circuits.
- o Speed-efficiency tradeoffs are such that an integrated logic circuit requires 20 per cent more parts than its discrete equivalent, and an integrated register circuit requires twice as many parts as its discrete counterpart.

Table II-8

COST VS. PROCUREMENT POLICY FOR A UNIT PARTS COMPLEMENT

Part	Quantity	Transistor Equivalent	Cost Extensions		
			As Is	Burn-in	Hi-Rel
Transistor	10,000	10,000	\$30,000	\$33,500	\$45,000
Diode	15,000	7,500	7,500	9,750	15,000
Comp. resistor	15,000	500	600	1,200	7,500 ^a
Film resistor	5,000	167	500	750	2,500
Mica capacitor	5,000	167	200	400	1,050
Total equivalent transistors		18,334			
Total parts cost			\$38,800	\$45,600	\$71,050

^aAssumes all film resistors used.

Table II-9

FAILURE RATE VS. PROCUREMENT POLICY FOR A
UNIT PARTS COMPLEMENT USING DISCRETE PARTS

Procurement Policy	Failure Rate %/1000 hr	Mean Time Between Failures (MTBF)
Computer grade	183	547
Computer grade with burn-in	92	1090
High reliability	55	1820

The breakdown of a unit system using integrated circuits is given in Table II-10.

Table II-10

DISTRIBUTION OF PARTS IN A UNIT SYSTEM

Part	Discrete	Transistor Equivalent	Uncorrected Integrated
Transistor	2,000	2,000	8,000
Diode	3,000	1,500	12,000
Comp. resistor	3,000	100	12,000
Film resistor	1,000	33	4,000
Mica capacitor	1,000	33	4,000
Total discrete transistor- equivalents		3,666	
Total uncorrected integrated equiv- alent parts			40,000

Using the assumptions given on p. 43, the number of integrated equivalent parts, corrected for speed and efficiency, is

$$2/3 \times 1.2 \times 40,000 + 1/3 \times 2 \times 40,000 = 58,666 .$$

Thus the required number of integrated circuits of complexity 10, 30, and 100 is

$$N_{10} = 5867$$

$$N_{30} = 1955$$

$$N_{100} = 587$$

and the failure rates for a unit parts complement computer, under the three procurement policies and using integrated circuits, are given in Table II-11.

Table II-12 summarizes the overall results. Before any conclusions are drawn, additional interpretation is required. Although the high-complexity, premium-priced (high-reliability), integrated-circuit approach seems to give the most system reliability per parts dollar, the quantities required for any really large program (e.g., data processing for a ballistic missile defense system) could easily challenge the capability of the industry to meet the need. For instance, if 100 systems, each of ten-unit-complement size, were built, 587,000 high-complexity integrated circuits would be required. It is doubtful that the stated cost-reliability relation could be maintained (pre-1968) with production requirements of this size.

Moreover, the above figures are sensitive to estimates of the efficiency-ratios of integrated vs. discrete circuits. If 1.5 and 3.0 are used instead of 1.2 and 2.0, the integrated approach appears less advantageous.

Also, consideration of the maintenance process and maintenance module size (Chap. VI) may show the integrated approach to be inefficient in maintained ground systems.

Table II-11
FAILURE RATE VS. PROCUREMENT POLICY FOR A
UNIT PARTS COMPLEMENT USING INTEGRATED CIRCUITS

Procurement Policy	10-Part Integrated Circuit		30-Part Integrated Circuit		100-Part Integrated Circuit	
	Failure Rate %/1000 hr	MTBF	Failure Rate %/1000 hr	MTBF	Failure Rate %/1000 hr	MTBF
Computer grade	154	650	116	863	78	1280
Computer grade with burn-in	106	944	67	1490	48	2080
High reliability	40	2500	29	3350	13	7700

Table II-12
SUMMARY OF COST-RELIABILITY RESULTS

Procurement Policy	MTBF (x unit complement, hr)				Cost (\$ per unit complement)			
	Discrete	IC10	IC30	IC100	Discrete	IC10	IC30	IC100
Computer grade	547	650	863	1280	38,800	66,460	46,860	31,260
Computer grade with burn-in	1090	944	1490	2080	45,600	79,420	56,120	37,320
High reliability	1820	2500	3350	7700	71,050	102,110	72,910	49,410

Further, entries in Table II-12 are parts costs only. The various costs and burdens imposed in the overall process of providing completed, installed, maintained systems may greatly decrease the significance of parts costs. The pro-rate cost of programming alone may approach the parts costs. Furthermore, the pro-rate costs of design for efficient utilization of high-complexity integrated circuits may be significant.

Finally, it may later be demonstrated that even the best MTBF given in Table II-12 is inadequate with realistic maintenance times, whereas the necessary step to some form of redundancy may make the simplest approach more than adequate.

The graphs in Figs. II-2 and II-3 show relationships among parts costs, system size, MTBF, and mechanization for non-redundant systems, with attempted indication of effects of some of the above factors.

7. PRESENT AND PREDICTED RELIABILITY LEVELS

Based on the available information, and assuming exponential failure distribution, Table II-13 gives estimates of present and predicted reliability levels for the part-types which have been evaluated. In spite of the fact that each manufacturer's sample that went into the compilation of Table II-13 claimed a confidence level of 90 per cent,[†] the spread of samples was so large that the accuracy of the entries is probably +100 per cent, -50 per cent.

Values in the table have 90-per-cent confidence levels, for 55°C maximum lead or case temperature, 70°C maximum

[†]The 90-per-cent confidence level samples were for 55°C maximum lead or case temperature, 70°C maximum hot-spot temperature, and 50-per-cent relative humidity.

Table II-13

PREDICTED PART FAILURE RATES
(%/1000 hr, 10-year average)

Part	Good 1965	Best 1965	Best 1968
Resistor (composition, metal film, tin oxide)	.0003	.0001	.00003
Capacitor (glass, mica)	.0003	.0001	.00003
Diode (silicon planar)	.005	.015	.0005
Transistor (silicon planar)	.01	.003	.001
Integrated circuit (silicon planar)			
10 equivalent parts	.02	.005	.0015
30 equivalent parts	.04	.009	.0025
100 equivalent parts	.07	.015	.0040

hot-spot temperature, and 50-per-cent maximum relative humidity.

Some discussion is in order on the integrated circuit entries, for two reasons. The first is the widely-publicized industry attitude that an integrated circuit can be just as reliable as a single transistor, because the manufacturing processes are identical. This just is not true; composite information from Borofsky [21] and failure mode/mechanism data (Table B-2) gives the approximate percentage contributions of planar process defects to part failure shown in Table II-14.

Table II-14

PER CENT CONTRIBUTION OF PLANAR PROCESS DEFECTS
TO PART FAILURE

Item	Defect	Percentage
1.	Package	17.0
2.	Gross (scratch, crack, foreign material, corrosive residue, etc.)	12.2
3.	Die bond to package	6.1
4.	Surface, contamination, passivation, diffusion	27.6
5.	Bonds to die	22.1
6.	Leads to terminals	7.2
7.	Deposited aluminum interconnections, windows, registration	6.4
8.	Silicon material	1.4

The first three items admittedly affect transistors and integrated circuitry equally. Items 5 and 6 are directly related to the number of external connections, which might be expected to increase somewhat with circuit complexity. A transistor has two internal leads (collector connected to case), a 10-equivalent-part circuit (e.g., 3 input NOR) has 6 leads, and 30-part and 100-part circuits are estimated to require 10 and 14 leads, respectively.

Items 4 and 8 are proportional to surface area. Although a 10-part circuit can be put on a chip not much larger than that required for a transistor (due to handling limitations), a 30-part circuit would require somewhat more

area, and a 100-part circuit considerably more at the present state of the mask fabrication and registration art.

Even with maximum topological cleverness, Item 7 is quite significant. It is estimated that the ratio of windows and interconnections to parts is 40 per cent for the 10-part circuit, 30 per cent for the 30-part circuit and 20 per cent for the circuit of 100 equivalent parts.

With the above estimates, it is possible to prepare a weighted relative failure rate figure for each complexity level. This is shown in Table II-15.

From the above very approximate analysis, it would appear that the failure rates for integrated circuits of 10, 30, and 100 equivalent parts might respectively be 2, 3.5, and 6 times that of a single transistor. Figures in the failure rate chart assume some improvement in these ratios in the future, and are scaled to the .005 entry for "best 1965, 10 equivalent parts," which essentially represents the most significant actual failure rate datum obtained [9].

With respect to external circuit failures induced by degradation of internal elements, integrated circuits appear to have a considerable statistical advantage, over discrete circuits. This advantage might be partially or completely lost, due to the compromises required in monolithic device design.

The second item for discussion related to integrated circuits is the rather poor overall showing of integrated circuits in the Compendium of Failure Statistics (Appendix E). At first, the figures seem to introduce serious doubts

Table II-15
RELATIVE FAILURE RATE FOR INTEGRATED CIRCUITS

Item	Weights				Weighted Percentages			
	Transistor	IC10	IC30	IC100	Transistor	IC10	IC30	IC100
1+2+3	1.0	1.0	1.0	1.0	35.3	35.3	35.3	35.3
5+6	1.0	3.0	5.0	7.0	29.3	87.9	145.5	205.1
4+8	1.0	2.0	4.0	8.0	29.0	58.0	116.0	232.0
7	1.0	4.0	9.0	20.0	6.5	25.6	57.6	128.0
Totals					100.0	206.8	355.4	600.4

about recommending their use in large-scale, ground-based computers where size and weight are not at a premium. As it is apparently possible to achieve .00002%/1000 hr failure rates with well-made solder joints or welds [9],[†] it is not reasonable to argue that interconnections will cause discrete circuits to be unreliable.

Moreover, it appears that the short test history of integrated circuits relative to discrete parts has not made it possible to prove their inherent reliability. The scatter chart of Fig. II-4 shows the relationship between the amount of testing and computed reliability. Results of 19 brute-force life tests were plotted, at the 90-percent confidence level, versus total unit hours of each test. The apparent correlation indicates the dependence of demonstrated reliability on test time. Indeed, the integrated circuits which established the point at 5×10^7 hr, .005%/1000 hr, may well possess the .0015 inherent reliability predicted for 1968, but some 90 million additional unit-hours would have to be accumulated for brute-force proof.

Reliability comparisons are further affected by the fact that an integrated circuit of ten equivalent parts will not directly replace ten discrete parts. At the present time, the discrete circuit will be somewhat more efficient in gain, bandwidth, noise rejection and time-race immunity than its integrated counterpart. Certain modified hybrid techniques, using miniature discrete transistors, overcome these obstacles at the expense of the as-yet-unknown reliability differential introduced by the assembly process.

[†]Also see Ref. 22, p. 53; Ref. 23, p. 159; Ref. 24, p. 211; and Ref. 25, p. 227.

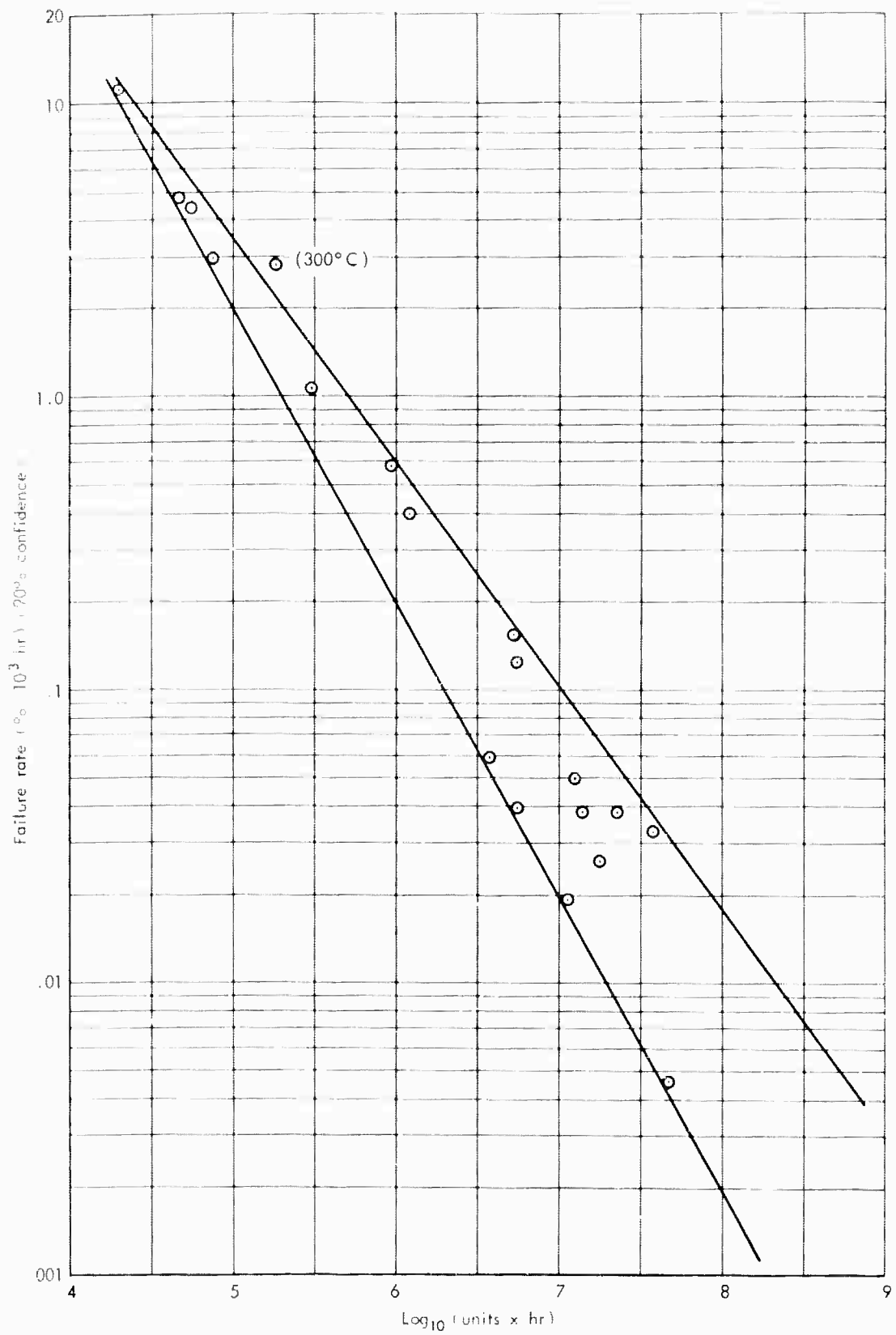


Fig. II-4 — Integrated circuit failure rate versus amount of testing

8. PREDICTED PART COSTS

The price of transistors, diodes, resistors, and capacitors is presently leveling off and, in some cases, showing a slight uptrend due in part to decreased demand. Improving processes and yield may decrease silicon planar diode and transistor prices.

Dramatic reductions are forecast in the integrated circuit area. Some reasonable assumptions must be made about the minimum cost of the package, packaging labor, and test labor or test equipment amortization. Table II-16 shows the predicted price chart reflecting a composite of industry predictions and minimum production cost estimates, even at high yield from the diffusion processes.

Table II-16

PREDICTED PRICE OF INTEGRATED CIRCUITS^a

Equivalent No. of Parts	Price, \$, 10,000 Quantities		
	1965	1968	1970
10 computer grade	10	5	2
10 high reliability	15	9	5
30 computer grade	20	10	4
30 high reliability	30	17	9
100 computer grade	40	20	8
100 high reliability	60	33	16

^aNote these are MIL computer grade components, not the so-called industrial-commercial grade units currently offered at a few dollars each.

Improvements in Materials, Processes, and Quality Control

The combined results of life-test and physics-of-failure studies are continuously being fed back into the manufacturing processes of competent suppliers. Needless to say, this is a diminishing-returns operation, at least where price is some sort of consideration.

As it appears that more parts-per-chip is a legitimate reliability objective for integrated circuits, improvements are required in mask-making and registration processes, and several seem to be on the way, as by-products of automatic machine-tool control, precision photolithography, and similar efforts.

Where the manufacturer does not have the impetus of a generously-funded, high-reliability program (Apollo, Polaris, Minuteman), his ambitions may be split between capturing some share of the coming industrial-commercial entertainment market in silicon devices and becoming one of a select group of qualified suppliers under high-reliability military specifications. Usually, qualification is at the expense of the buyer, and in some cases is quite expensive. Included in the brute-force life test requirements of MIL-R-38100A, for instance, is the following case:

Class Z, .0001%/1000 hr, 90 per cent confidence, requires testing of 23,026,000 parts for 1000 hr with no failures for qualification.

A user requiring a million half-watt composition resistors qualified as above, would have to pick up the "reliability overhead" of 23 million additional resistors, a test rack consuming 11.5 megawatts of power, and some six man-years of before-and-after measurement (at one second per resistor).

Similar fascinating requirements are found scattered profusely throughout the many high-reliability military specifications. It is exactly this situation which lends support to the physics-of-failure, quality-control oriented approach to lower-cost, higher-inherent-reliability parts production.

Possible Breakthroughs

Reduction to practice of thin-film active element production, or demonstration of high-reliability in hybrid methods using discrete transistors could introduce an important alternative approach to high-availability ground-based computer production. The significant potential of the hybrid is in the increased circuit efficiency due to more stable resistor and capacitor values, freedom from internal compromise, and higher net yield, as passive and active components may be tested (and burned-in) separately before assembly. The IBM System/360 is committed to a hybrid (flip-chip) approach, and large quantity quotations of \$1.00-2.50 per flip-flop have been obtained from hybrid suppliers for other systems.

Other possible approaches under development are metal-oxide-silicon integrated circuits and field-effect transistor circuits. It is doubtful that high-reliability procurements initiated in 1965 should depend on the to-be-demonstrated reliability of these devices

9. CONCLUSIONS AND RECOMMENDATIONS ON PARTS

Semiconductor Devices

Transistors and diodes have three conveniently related characteristics:

- o They account for most of the parts cost (84-97 per cent in the examples above);
- o They are responsible for most of the failures (viz., 95 per cent);
- o Reliability improvement by conventional means costs the least percentage-wise (50-100 per cent for transistors and diodes; 400-900 per cent for resistors and capacitors).

Furthermore, the cost of parts is only a moderate fraction of the cost of checked out, delivered, installed computer systems complete with working programs. Therefore, the following general requirements should be imposed on semiconductor procurement:

- o Select potential suppliers on the basis of proved reputation in, and apparent large-scale commitment to, the high-reliability market;
- o Require a two-stage 100 per cent screening process on parts, plus mechanical and environmental testing on a sampling basis, such as specified in MIL-S-19500D. The overall process should be as shown in the flowchart given in Fig. II-5.

The second-stage burn-in and test may be performed by either the supplier or the user. The decision may well be an economic one, although availability of results of second-stage burn-in is an invaluable aid in selecting a vendor. Obviously, any failures should be treated as rare

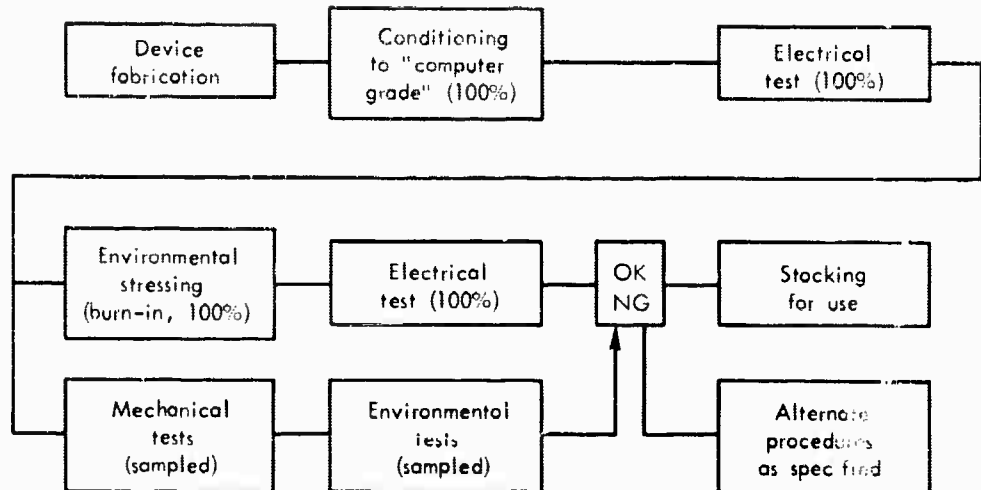


Fig. II-5—Test process for semiconductor device procurement

and valuable items and subjected to a carefully-designed post-mortem examination, usually by the device manufacturer.

Require or perform some form of long-term life testing on samples of devices taken after second-stage screening, to

- o Detect any changes in quality not caught by screening and sampling tests;
- o Uncover any new long-term inherent failure mechanisms;
- o Be sure that burn-in is not accelerating any long-term failure mechanisms.

All of these recommendations apply equally to integrated circuit procurement. Although, on a "parts-only" basis, it appears logical to insist on high-complexity integrated circuit implementation, such a decision should be subject to evaluation of the following:

- o Effect on maintenance module size and cost, the number and type of spares to be stocked, and the equipment and training required for depot maintenance;
- o Effect on design cost and complexity and interconnection reliability due to topological restrictions and high interconnection density imposed by small size;
- o Actual price, availability and net (post-screening) yield of high-complexity circuits at time of system procurement.

For procurements initiated in calendar 1965, it appears that reputable digital-systems manufacturers could present valid cases for either discrete or monolithic silicon integrated-circuit implementation. Also, it is possible that manufacturers committed to, and skilled in, the modified hybrid ("flip-chip," etc.) approach could compete if some objective demonstration of the reliability of the configuration can be presented.

Some recently published information [26,27] indicates that the "best 1968" reliability levels for transistors and diodes are achievable today, in organizations having complete control, from the raw material procurement to equipment installation and maintenance. It is interesting to note that the system, Bell Electronic Switching System No. 1 [8], is duplexed for reliability at the systems level, has a target of 40 years mean time to dual failure

with one- to three-hour service time, and is constructed with discrete parts on etched cards and wire-wrapped backboards.

Dickenson [28] describes a hybrid thin-film approach for a high-reliability space mission using triple voting redundancy. Logic modules of complexity up to 14 parts are used with an estimated individual failure rate of .03 to .04%/1000 hr.

Fagg, et al., [29] and Davis, et al., [30] describe a series of commercial computer implementations using essentially the same hybrid approach.

Resistors

As costs and reliability effects are negligible, metal-film or tin oxide resistors should be used wherever the narrowed degradation tolerances, and availability of a larger selection of nominal values, permit higher circuit efficiency. A decision to use composition resistors elsewhere appears optimal, unless there is some effect on spare parts costs. For all resistors, temperature-cycling followed by 100 per cent inspection on a limit bridge is recommended to catch gross genetic defectives.

Capacitors

Conventional dipped-mica capacitors may be used, with 100 per cent voltage-temperature burn-in and limit-bridging. High-reliability units might alternatively be selected, but there seems to be more assurance (if only emotional) in 100 per cent inspection at the user's site, and the latter

may cost considerably less. Mylar capacitors should be handled as above.

In-Plant Parts Handling

Parts of initially excellent quality are exposed to a high risk of traumatic experience along the path from receiving inspection to equipment installation. Innumerable examples of mistreatment, from a single resistor to entire computer-logic sections, have been discovered. A typical quote is "Oh, yeah, old serial number 3 is always a pain to keep on the air. That's because, when we were checking it out, a regulator went out in the lab power supply and the voltage went from 12 up to 50. There's still a spot on the ceiling where one of the electrolytics exploded." But old number 3 was nevertheless shipped.

The following recommendations should be imposed on all system suppliers--and, in fact, any not already incorporating most of them (particularly in integrated circuit work) should be viewed with suspicion.

- o After receiving inspection, store parts methodically in a known environment and in a manner which prevents physical damage. Preferably, parts are stored in the carriers in which they will later be delivered to the assemblers.
- o Assemble under scrupulously controlled conditions of environment, cleanliness, and operator aptitudes, training, and integrity.
- o Devise test and checkout equipment and its interconnections such that there is a vanishing probability of part overstress through misuse or malfunction. Conduct suitable selection and indoctrination of test and checkout personnel to minimize overstress but to insure it is reported if it occurs.

Circuit and system design will, of course, incorporate adequate part derating factors in normal operation. Design must also minimize probability of overstress in all failure modes. Power supply regulator shorts are probably the outstanding offenders, but every system contains numerous disastrous possibilities which should be evaluated. Protection against transients from power lines, during normal and abnormal power turn-on and shut-down, must of course be provided.

If high pointwise availability is the goal, then, despite the possible two-to-one reliability differential, it is better to have the equipment on continuously than to turn it off and on as few as two times daily, assuming there is no possibility of power transients. This is principally true because of the ever-present chance of over-voltage surges during the transient period and the impossibility of performing error checks when the machine is inoperative. If, on the other hand, the goal is to maximize the "interval availability," namely, the fraction of a specified interval of time that the machine is operating, then there might be a stronger case for turning the machine off when it is not needed.

Reliability of the operating environment is essential. Humidity and temperature control should have reliability comparable to that of the computer, and of course must be interlocked.

With all the above precautions, the largest single hazard to parts has yet to be mentioned: the maintenance man. Nearly every authority, at some stage of the discussion introduces a statement such as "and once you get

it going, leave it alone. Don't let anybody into it with clip leads or probes (etc.)." Assuming good design for maintenance and good diagnostic equipment and procedures, the best insurance for parts is to make sure that every maintenance technician or engineer avoids unauthorized procedures and reports every "accident" faithfully.

A final note on parts: All comments on screening, burn-in, handling, and storage should apply, where feasible, equally to spares.

Chapter III

THE RELIABLE COMPUTER

1. INTRODUCTION

This chapter discusses the problems of designing and building a single, reliable computer. The possibility of achieving reliability by using more than one computer is deferred until Chap. IV, but use of internal redundancy at various levels is discussed here.[†]

The major influences on reliability are not the schemes, such as redundancy, which are sometimes implemented, but the fundamental methods by which reliable parts are converted into an operating computer--for it certainly isn't true that reliable parts insure a reliable machine. For this reason, considerable space goes to a philosophy of good circuit design.

The process of logical design also requires attention--errors in logical design or construction can show up well after the computer is in the field, with serious and usually irreversible consequences for military operations.

Much has been written in recent years about error detection, and the subject is an important one. We have chosen to expand the topic under the heading "failure detection" wherein "errors" are part malfunctions which do not produce immediate signal errors, as well as the usual

[†]The formal analysis of the redundant computer is performed in Appendix A, and only selected results will be given here.

signal errors themselves. Familiar terms such as "error," "detection," "correction," etc., are carefully defined, and methods of coping with these failures are discussed.

This chapter concerns only those circuits which protect against failures. Chapter VI discusses fault diagnosis, isolation, and correction under program control. Such programs, as will be seen, can powerfully affect repair time and, hence, availability.

2. CIRCUIT DESIGN

As discussed in Chap. II, any part "failure" definition other than passage to a limit value is arbitrary and should be related to the intended application. Practical circuits must be designed with some allowance for part parameter variation with environment and time, and some initial parameter range due to manufacturing variations.

Given a circuit configuration (schematic diagram) and perfect knowledge of degradation behavior of the parts, it is theoretically possible to determine circuit reliability-versus-performance relationships for some stated design life. For any but the most elementary configurations, the mathematical task involved in determining such relations becomes impressively complex.

The minimum circuit of interest contains three or four inputs, one or two outputs, two or three supply voltages, and ten to twenty parts. A first-order characterization of a resistor or capacitor requires but a single parameter. For a transistor or diode, even the simplest mathematical model requires several parameters and some mode of approximating nonlinear behavior in various operating regions.

The engineer with a slide rule attempting to analyze a circuit must be content with crude piecewise-linear approximations and one-pole models of the active elements, such as those of Ebers and Moll [1] or Beaufoy and Sparkes [2]. The accuracy of this kind of work (relation to observed circuit behavior) is ± 10 per cent at best, and -90, +500 per cent at worst, depending on circuit speed and complexity.

Developments in the last decade dramatize the need for "better tools", and the attempts to provide computational relief and more accurate modeling show varying degrees of utility and success [3-8]. With integrated circuits, sophisticated models and high-powered computation are needed, because of

- o The inherently more "distributed" nature of integrated circuits;
- o The near-impossibility of "breadboarding" and "laboratory design" by successive modifications.

It should be emphasized that the above work relates to analysis only. True circuit synthesis by automated means is, at best, still in the "laboratory curiosity" stage.

The usual process of circuit synthesis proceeds somewhat as follows:

- o End-to-end system specifications are somehow subdivided into circuit functions by cooperative action of systems, logic, and circuit designers;
- o The circuit designer takes a functional specification for a circuit and--through some combination of experience, intuition, and plagiarism (research)--selects one or more tentative configurations and sets of active elements;

- o Using first-order (or zero-order, e.g., "educated guess") approximations, a preliminary set of part values is computed;
- o The preliminary configuration is analyzed, using better models and approximations, if possible;
- o Shortcomings are corrected by parameter change or substitution of better parts, as required;
- o The iterative process of analyze-modify is continued until performance is satisfactory or hope is abandoned and a new configuration is sought.

The necessary implements of the process are models and computation. Models may be 1) lumped-constant equivalent circuits, valid over certain specific parts of the operating region, 2) various degrees of nonlinear, distributed-constant, or 3) true mathematical analog representations. Computational aids fall into two general categories: analysis and simulation. Analysis, of course, is ideal up to the point at which complexity renders it impractical. Simulation, both analog and digital, has been utilized, with incremental-digital approaches undergoing considerable investigation at this time.

Bogey Design

All circuit design systems require some guiding philosophy which defines theoretical circuit failure as a function of theoretical part behavior. The earliest and simplest philosophy, regrettably still followed in some organizations, is "bogey" design.

In its fundamental form, bogey design assumes that all part parameters are at their new-nominal (as-labeled)

values, and will stay there forever. The only reasonable defense that might be used by a thinking bogey-designer is that positive and negative deviations are equally likely, small deviations are more likely than large, and the effect of many circuit topologies then gives some sort of statistical protection. A slightly advanced form of bogey design admits part variation over the manufacturing tolerance range, but no more. The effect of part selection by the supplier and the resultant rectangular distribution in as-delivered parameters, is obvious.

Worst-Case Design

The wave of reaction to bogey design, started by the obvious fall-off of reliability as systems become more complex, led to the formation of various "worst-case" design philosophies. The purest and simplest of these (called "worst-worst-case" in some circles) operates as follows:

Absolute end-of-life limits are provided for all part parameters. The designer does not consider whether these are really absolute or actually some sort of statistical limit, but designs the circuit so that it meets specifications with all parameter values simultaneously at their most unfavorable limits. This usually requires selecting a different set of limits for analysis with respect to each operating specification.

The beauty of this process is its mathematical simplicity. Sets of limits are selected, usually by inspection, or at worst by a few coarse trials; the usually few part parameters contributing to the function under evaluation are then determined by solution of sets of simultaneous inequalities. If the limits are properly selected and if there are no mathematical errors, the circuits always work, and continue to work.

The major disadvantage here is the inefficiency of over-design. Ultra-conservative worst-case protection against degradation increases the number of parts in a given system and the vulnerability to part failure of the "limit" sort (open, short, anomalous degradation).

A reasonable compromise, now practiced by many responsible organizations delivering excellent equipment, is worst-case design with narrowed tolerances. They use the same mathematical procedures but they perform limit selection on a statistical basis to yield narrower spreads.

Another approach is the so-called "Taylor worst-case" design, wherein the specifications must be met when any one part is at its most unfavorable extreme, with all others either at

- o New-nominal, or
- o The most unfavorable extremes of manufacturing tolerance.

This technique is not often used because the mathematical labor is excessive, and the compromise varies with circuit configuration.

Statistical Design

Statistical design is the theoretically-ideal design philosophy because it considers the behavior of the circuit throughout the range of parameter variation, instead of just at the limits of parameter values. The required information is

- o Distributions of part parameters with stress and time;
- o Limits of stress and the design life;
- o Required reliability of the individual circuit, i.e., probability of meeting specifications in the environment, over the design life.

The design process, then, somehow yields reliable parameter values or indicates that the requirement cannot be met with the given parts and configuration. Ideally, distributions of stresses would be given rather than limits, and, as ranges of parameters would result in many cases, an additional constraint might be introduced, such as minimum power consumption. Obviously an engineer with a slide rule cannot perform statistical design. Attempts at aids to statistical analysis have produced isolated results in man-computer circuit synthesis [5,9]. But at present sheer complexity prevents the large-scale use of statistical design.

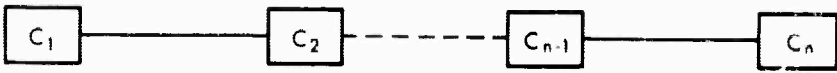
Protection Against Part Failure

Design to protect against part failure, rather than degradation, may exist at any level from circuit through module, or subassembly to the entire system. The philosophy is generally referred to as redundancy at the selected level. Two basic approaches to redundancy may be classified as switching and paralleling. If the non-redundant system is considered as a series chain of elements (subsystems, circuits, parts) which fails when any element fails, it may be represented as shown in Fig. III-1a.

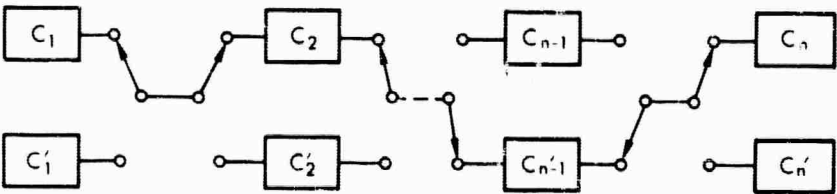
Switching redundancy, in its simplest form, implies the existence of a spare for each element and provision for switching it into the system, as shown in Fig. III-1b.

There are two additional requirements:

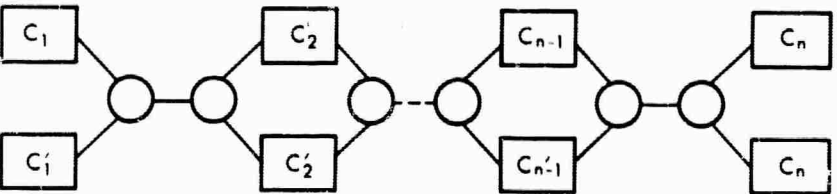
- o Means--either human or automatic--for detecting failure and localizing to a particular element;



a. Non-redundant series chain circuit



b. Switching redundancy



c. Parallel redundancy

Fig. III-1 —Redundant circuits

- o Design the system and the switch such that time required to switch does not affect system operation.

If several of the elements are identical, only one spare of that type need be available, provided the switch is suitable.

Figure III-1c shows parallel redundancy. Here the circles represent interconnections designed such that failure of an element in any possible mode cannot prevent the associated parallel element from performing its function.

Switching redundancy is sometimes referred to as "standby" redundancy. The spares may or may not be energized prior to use. For computing reliability of elements subject to wearout failure under power (light bulbs, electro-mechanical devices) this distinction is extremely significant. For systems which are predominantly semiconductors, the distinction is insignificant, considering the cost of the added switching complexity to apply power to a spare.

At the part level, switching redundancy is impractical, as the reliability of a switch is, at best, similar to that of a part. Parallel redundancy at the part level is achieved by some sort of series-parallel connection of single parts or "subcircuits" consisting of a few parts each. Further discussion on circuit redundancy is deferred to Sec. III-4.

Recommendations in the Procurement of Good Circuit Design

Several requirements for responsible, reliable circuit design are listed below. These might be considered as checkpoints on the integrity of design.

- o There must be a definitely stated, defensible design philosophy.
- o Rigorous adherence to the philosophy is essential, and high-level technical-administrative approval should be required for allowance of any exceptions or discrepancies.

Part-types must be rationally selected.

- o Consider present and predicted availability when establishing design tolerances.
- o Require considerable investigation before selecting, and in determining limits for, "novel" parts. This might include four-layer devices, unijunction transistors, thermistors, field-effect transistors, tunnel diodes, and multiaperture ferrite devices.
- o Avoid specification of low-yield parts which may become unavailable due to slight shifts in material and process parameters (e.g., extremely high-gain transistors).
- o Avoid designs specifying parts requiring selection in narrow ranges or close matching and tracking.
- o Avoid dependence on uncontrolled part parameters for circuit function (e.g., reverse recovery time of coupling diodes for transistor turnoff).

Intensive effort must be devoted to anticipation and evaluation of system-imposed requirements on circuit specifications. Some examples follow:

- o Provision for guarding against noise;
- o Determination of true nature of interface signals, rather than acceptance of superficial specifications;
- o Regardless of the design philosophy, recognizing situations where worst-case reasoning must be applied (as in output signals and noise from core memory stacks, where any core may be interrogated with any combination of overall information pattern and previous history).

In human-slide rule systems, prefer worst-case design with intelligent compromises in the selection of end-of-life tolerance limits and mathematical models of the components.

In human-computer systems, evaluate the extent and nature of computer aids. Beware of systems under development; insist on acceptable proof of operability of models and computing techniques.

In any case, look for a stated design philosophy, organized review procedures, and well-kept engineering notebooks. Regardless of the excellence of the above procedures, get results of lab verification of single circuits and circuits operated in all reasonably achievable combinations of loading, noise, layout, interconnection, and stress.

3. LOGICAL DESIGN

The logical designer interacts with the system designer and the circuit designer to generate an interconnection of logic elements which will

- o Functionally implement the system specifications, and
- o Operate electrically when interconnected, without violating established circuit limitations.

The logic designer's responsibility varies widely from organization to organization. He may work directly from system specifications, thus annexing the system design function. He may generate actual diode network schematics, thus encroaching on circuit design. Usually, though, the logic designer writes a set of logical equations, or generates a logic schematic diagram, subject to the constraints of system design and circuit interconnection complexity limits. In the most sophisticated systems, the logic designer (after doing his own private scratch-work in his favorite form) makes direct symbolic entries, representing his equations, on a suitably human-engineered keypunch form. Some subset of the following steps is then executed, depending on the magnitude of the system.

- 1) Keypunch original or modifications and verify.
- 2) Run computer routine to check for violation of circuit limitations, and various consistency checks.
- 3) Cycle 1 and 2 until deck passes.
- 4) Add this segment to logic simulator.
- 5) Run logic simulator--results go to logic and system designer for approval.
- 6) Cycle 1-5 until buildable portion is complete.
- 7) Perform layout according to minimum-wire-length, noise, and other rules. Check dynamic loading, and list discrepancies for logic, system, or circuit designer action.

- 8) Generate and check control tape for production of backboard by automatic wire wrap or automatic multilayer laminate method.
- 9) Produce by-products:
 - a) Bill of materials
 - b) Layout diagrams
 - c) Logical block diagrams
 - d) Maintenance manual.

An adequate present-day system would include steps 1 and 2, with the logic designer or a specialist producing the layout and the computer producing wiring or interconnection instructions. The complete system described obviously provides maximum protection against human transcription errors, wiring errors, and design errors with respect to system and circuits.

The value of added steps is difficult to estimate. If designed specially for a production run of only a few systems, the cost of logic simulation and optimum-layout programs would be prohibitive. However, organizations already possessing such programs may be writing them off across large-volume production, with obvious gains in value.

A complete design automation procedure has peripheral value, in that programs may be run on a simulator (at a price) in advance of system construction, and maintenance and checkout procedures may likewise be established. An estimated 80 per cent of the design-production errors in all-human systems are the result of clerical and technical mistakes, rather than conceptual design flaws. Design and production automation reduces these.

Probably the first delivered unit of every large-scale computer contains several logical, wiring, circuit-design, or operating-program errors. The normal process of checkout, delivery, and customer feedback on high-production (100 or more) computers results in detection and elimination by retrofit and field modification of all, or very nearly all, such errors.

Such errors could, of course, be eliminated at the source by exercising each computer by another (presumably perfect) computer which simulates all possible modes of operation. At any predictable computation rate, the time required would, of course, be astronomical. However, should a dual- or multi-computer approach be selected, it is feasible to consider the "exhaustive exercise" as one check to be carried out, over the years, at all installations. Any tentative error which is reported could be verified at a second location, before implementing corrective procedures. The cost involves writing and checking the simulator program, which could nearly double the total programming cost.

4. CIRCUIT REDUNDANCY

A circuit is said to be n -fold redundant if that circuit is replicated n times and the output of the aggregate is taken to be the output which is produced by the majority of the circuits. This is schematically illustrated in Fig. III-2.

To prevent a tie, n is assumed always to be odd. A "circuit" may, in theory, consist of a single part or an entire computer.

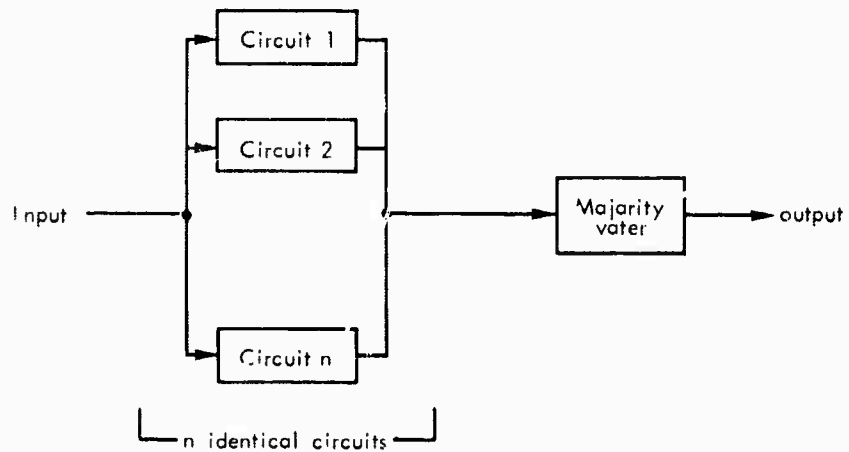


Fig. III-2—Redundant circuits with majority voter

In an n -fold redundant subsystem, no service is performed until the entire redundant subsystem fails, i.e., when more than half of the circuits have failed. At that time, service is performed and all n circuits are restored to an operative condition.

Only results pertaining to the use of redundancy will be presented here. Appendix A presents a complete discussion and derivation of these results.

Let $P(t)$ be the probability that the computer is on at time t . First, consider $P(t)$ as a function of computer size and individual part failure rate: In Sec. II-7, the unit part complement was defined and taken to be 18,334 equivalent transistors. A unit part complement of this size implies the failure rates (for present-day parts) shown in Table II-9. In addition to these failure rates, others derive from either using the future predictions of Table II-13, or assuming more unit part complements per computing system.

In general, when service is possible, the transient phase is not important and the asymptotic probability, P_{∞} , is the important measure of availability.[†]

Let μ = the service rate [$1/\mu$ = the mean time to repair (MTTR)], and $N\lambda$ = total number of parts times the part failure rate. $N = 18,334k$, where k is the number of unit part complements per system. Then P_{∞} as a function of $N\lambda$ for the non-redundant computer is shown in Fig. A-12 for various service rates. With n -fold redundancy, let M be the level of the redundancy, i.e., M gives the number of n -fold redundant modules in the computer. Figures A-16 to A-19 show P_{∞} for a system having three-fold redundancy and $M = 1, 10, 10^2, 10^3$. For the same set of μ 's and M 's, the availability of a five-fold redundant computer is shown in Figs. A-20 to A-23.

A few of these results are collected together and shown in Figs. III-3 and III-4 with some pertinent annotations on present and predicted performance.

[†] $P_{\infty} = \lim_{t \rightarrow \infty} P(t)$. See Appendix A for a discussion of this point.

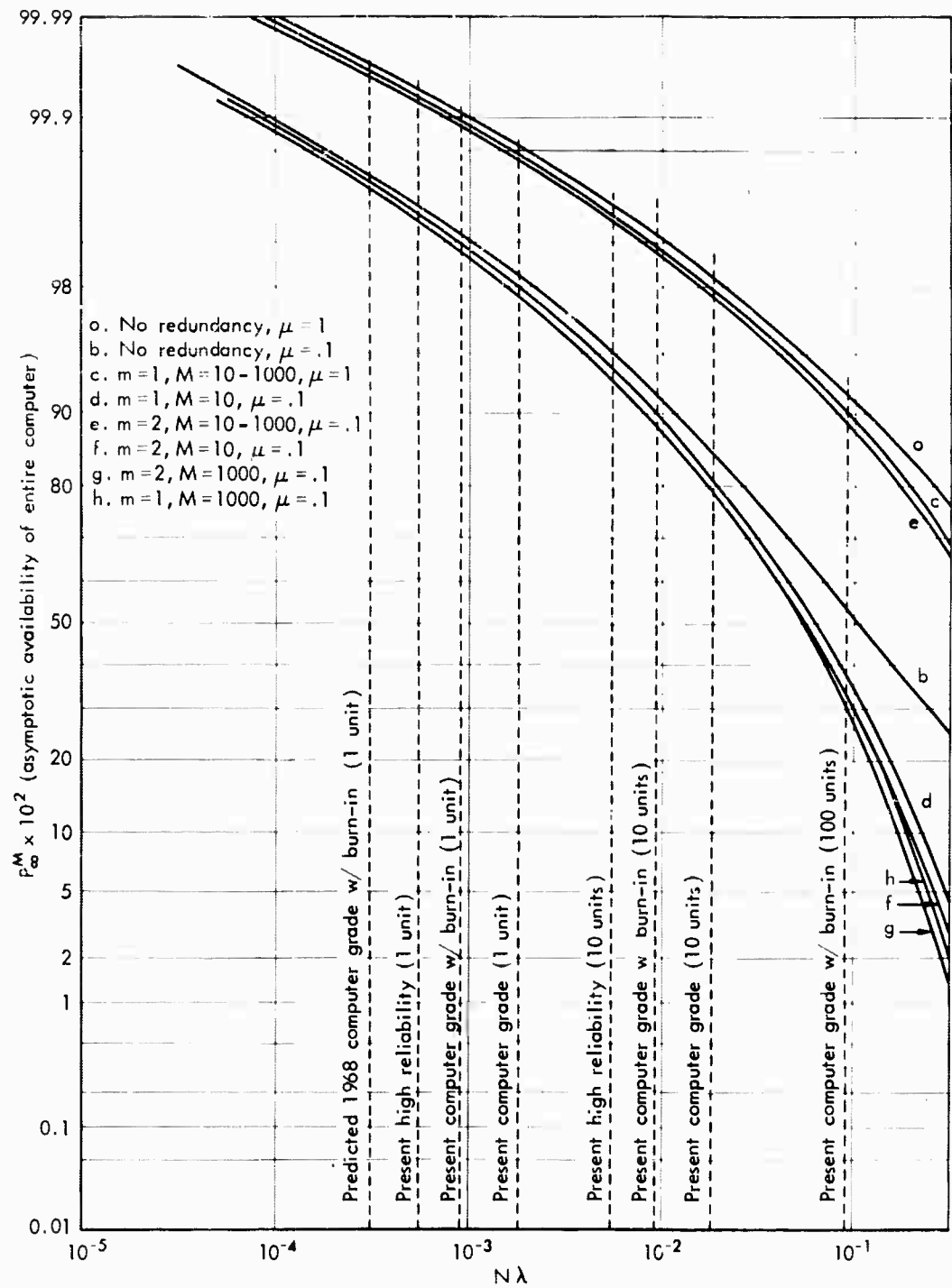


Fig. III-3 — Asymptotic availability of redundant computers (exponential service)

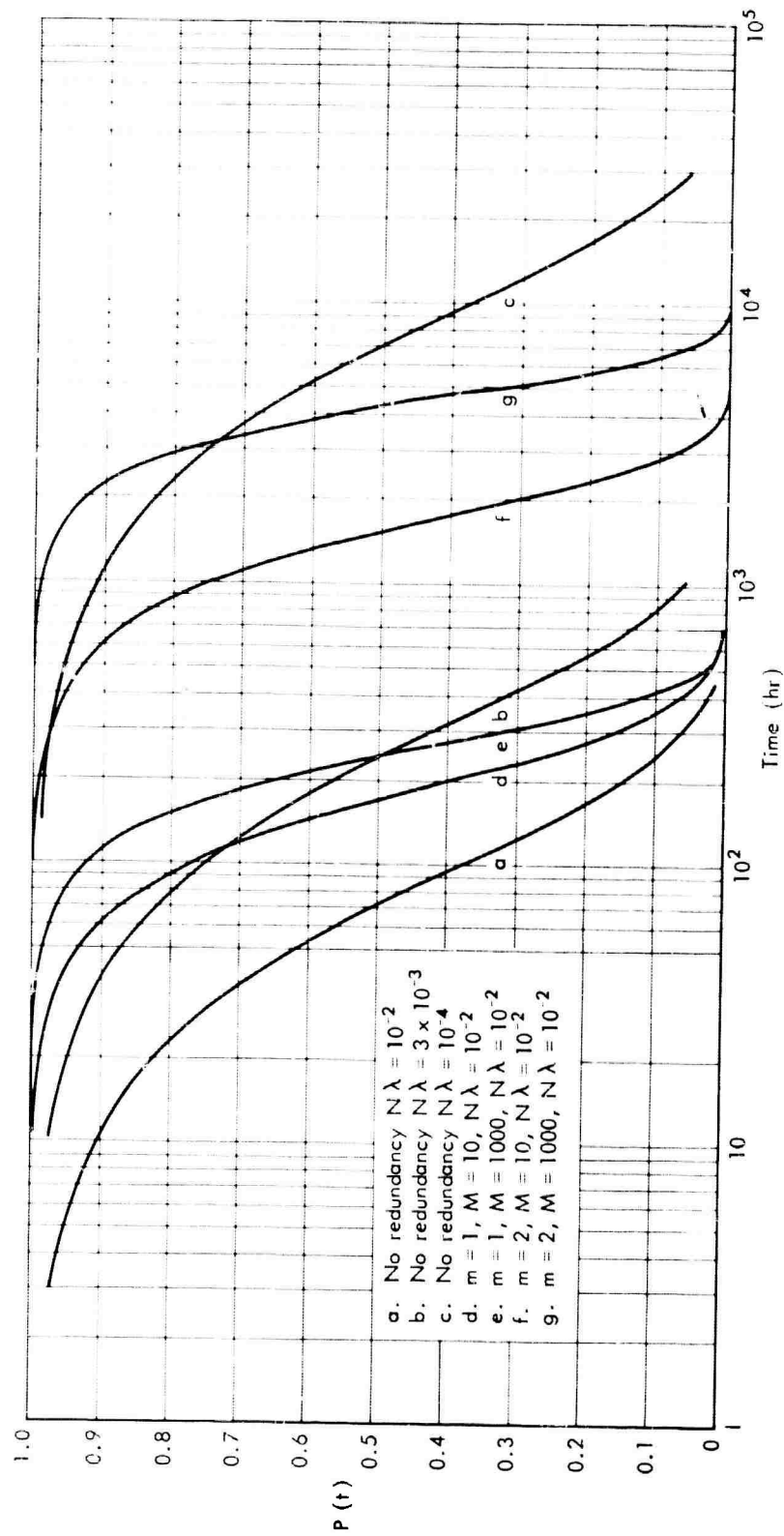


Fig. II-4—Availability of redundant computers
(no service)

The transient phase for redundant computers is shown here only for the no-service case; when service is available, there are better methods than redundancy, as will be shown in Chap. IV.

All the results presented in Appendix A will not be duplicated here. The reader who wishes to pursue further the subject of availability of a redundant computer should find the values of failure rate and number of parts from Chap. II (or use his own values), choose a service rate μ , and a level of redundancy M , then use the graphs in Appendix A to find either P_{∞} (or $P(t)$) for some special cases.

5. FAILURE DETECTION

Clarifying terminology at the outset requires the following definitions and discussion:

Defect--An inadequacy in the logic, wiring, or program of the machine as built. Some defects may indeed be detectable by the procedures to be discussed. Obvious cases will be noted, but the possible existence of defects will in general be ignored. (Section III-3 discusses prevention of logic and wiring defects, and Sec. V-1 discusses programming-coding defects.) Installation of an initially-failed part would also constitute a defect, but exhaustive module tests before assembly are assumed to preclude this possibility.

Maintenance Modules--The set of smallest field-replaceable elements. If a connector or an inter-connection fails (hopefully an extremely rare event),

it is probably repaired in the field, as it seems unreasonable to define the entire backboard as a maintenance module.

Fault--Abnormal electrical behavior of a maintenance module. A fault may be the result of part failure or unpredicted stress. A permanent fault results from part failure or (improbably) from a permanent occurrence of unpredicted stress. A transient fault usually is the result of occurrence of an unpredicted power, noise, or signal condition of short duration; a less likely cause is occurrence of a short-term unpredicted stress combination (temperature surge, physical shock); a possible (though improbable) cause is the excursion of a part parameter out of, then back within, degradation limits.

Error--An instance of incorrect functional performance by the system (e.g., least significant bit of accumulator is always "one"; wrong result of division of 32,169 by 20,447; launched 39 interceptors at a low-flying ptarmigan). An error is reproducible if it always occurs when the proper circumstances are established within the machine; an error is transient if it occurs once and efforts to re-induce it fail. Note that a transient fault, or a permanent fault which is located and corrected, may or may not have produced an error. Also, transient errors are associated with transient faults, and reproducible errors with permanent faults. As subsequently used, "fault" alone refers to permanent faults.

Detection--The process of determining that a fault or error has occurred. Also called recognition.

Correction--Applied to an error, removal of the error before the erroneous information is used.

Location--The process of determining which module is faulty, or which part has failed. Also called isolation.

Repair--Replacement of a faulty module or a failed part.

Service--The sequential execution of all of the above five processes.

There are three levels at which the above five functions may be performed (distinguishing fault detection from error detection):

- o The human level, presumed self-explanatory in all five cases.
- o The machine execution level--that is, in the process of executing a stored program. (Also called the "program" or "software" level.)
- o The machine implementation level--that is, in the design and physical construction (also called "built-in," "automatic," or "hardware" level).

For convenience, the adjectives "programmed" and "built-in" refer here to the machine execution and machine implementation levels, respectively. What follows next is a detailed examination of the various possible models of operating, excluding the human level.

With respect to repair, "programmed repair" is, of course, impossible; i.e., there is no executable computer code which will cause the computer to repair itself. "Built-in repair" is essentially synonymous with switching redundancy (cf. p. 72).

"Error detection and correction" permits a system to operate even after faults occur. Literature on the logic and methods of detection and correction is extensive [10-14]. It is essential to determine the value of this form of protection relative to the system application. The value, in turn, is a function of the reliability increase, the effects on operating speed, and the equipment cost of

- o The amount of detection-correction provided and
- o The effect of partitioning between programmed and built-in methods.

Errors are found in information; faults are found in equipment. Errors may be broadly classified as follows:

Transfer errors--Information which should be identical at two points separated in space-time is in fact different. Memory transfer errors occur in reading from storage, or less often on writing into storage. Input-output transfer errors occur in many forms, depending on the actual equipment and transmission system.

Operational errors--The result of some arithmetic or logical operation is incorrect.

Control errors--That portion of the machine which identifies and sequences operations performed improperly (e.g., subtraction is performed instead of addition; a branch is executed to the wrong instruction).

The value of detecting and correcting an error depends on the probability of its occurrence, the cost of protection, and the consequence of letting the error exist. Some

systems (ballistic missile defense) are unusual in that they are subjected to extremely short periods of peak demand at extremely infrequent intervals--perhaps one peak demand, or no peak demands at all, may occur over the entire design life. The consequences of an error during a peak demand period, however, are extremely severe.

A very simple model may give some insight into the values involved. Assume peak demands of up to 100 seconds duration, occurring once in ten years. Assume system downtimes of 1, 10, or 100 hr per year. (For ten-hour repair time this corresponds to MTBF of 87,600, 8760, and 876 hr, respectively). Further assume transient error rates of 10 or 100 per year (36 days and 3.6 days mean error-free time, respectively). Figure III-5 shows the relative risk as a function of demand duration, for various combinations of downtime and error rate. Risk was computed assuming transient error duration is negligible relative to demand duration, and uptime periods are large relative to demand duration, giving

$$\text{Risk} = \frac{ET + 3600D}{10 \times 365 \times 24 \times 3600} ,$$

where E = number of errors per year, D = number of downtime hours per year, and T = duration of peak demand in seconds.

The curves indicate that, if risk values of the order of 10^{-3} are adequate, transient error effects are negligible. But with risk values approaching 10^{-5} , transient errors are relatively significant. Transient error rates are difficult to predict. It is feasible and mandatory

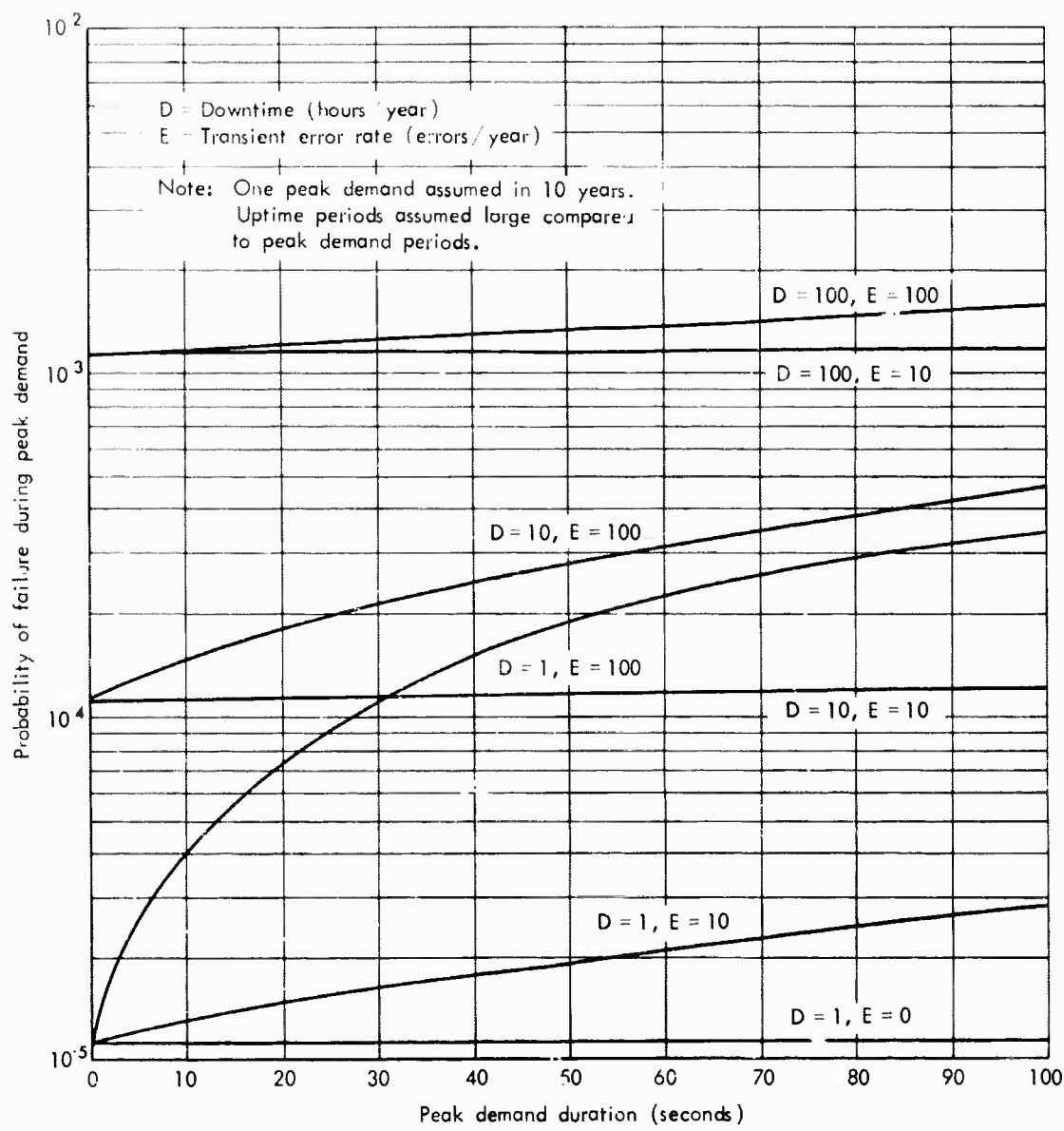


Fig. III-5—Probability of failure during peak demand period

to provide all reasonable system and circuit protections against transient errors resulting from power line behavior, electromagnetic interference, and the devices involved in information transfer.

The remaining actions that may be taken are error detection and logging, and error correction. Detecting transfer errors is a relatively simple process, usually accomplished by parity checking. Control and operational errors are not so easily detected, because the nature of decoding, arithmetic, and logical operations is such as to destroy simple internal relations, such as parity. Fortuitous exceptions are numerous and it is reasonable to consider error detection in these cases, where the added cost is justified. Examples are:

- o Forbidden operation code detection;
- o Forbidden digit checks;
- o One-only checks on decoding matrix outputs.

Redundancy in arithmetic, logical, and control equipment, of course, introduces error protection, but at a parts cost at least proportional to the order of redundancy.

Built-in parity checking requires addition of one bit to all transfer paths and storage locations. Also, the added circuits to generate and check the parity are required. For a reasonably large machine (3000 flip-flops, 25,000 gates, 36-bit parallel operations) parity checking and generation should not add more than 3 per cent to the parts complement.[†] Programmed parity checking is possible,

[†]Methods of implementing parity checks are discussed in Refs. [12-14].

but still effectively requires the extra bit of storage in that one bit must be redundantly "wasted" in each word which is checked. Also, this is a time-consuming process which might introduce a hidden cost in the form of increased overall speed requirement.

In a single, series-chain (non-redundant) system, the utility of single-transfer error detection differs markedly for destructive and non-destructive transfer. In destructive transfer, information in a register (usually a set of magnetic cores) transfers to a second register, and the first register simultaneously clears. If an error occurs in the transfer, the information is irrevocably incorrect, even though the error was transient. In non-destructive transfer, the content of the first (sending) register remains, at least until the receiving register completes error checking. In this case re-transmission can correct a transfer error.

In dual systems that operate in synchronism, detection of an error by one machine can be used to switch the non-erratic machine on-line. Multiple systems provide many other alternatives.

The correction of single transfer errors by built-in means requires the addition of considerable equipment. Six additional bits must be generated and carried for correction of a single error in a 36-bit data word. Depending on the number of transfer paths accommodated, up to 50 per cent additional circuits might be required.

For a single non-redundant computer of about 100,000 equivalent-transistor complexity, it is unlikely that MTBF of more than 1000 hr could be achieved. For repair

times longer than one hour, the down time component of risk would be greater than 10^{-3} . This apparently makes error detection and correction unnecessary. If dual- or multi-computer redundancy is provided, error correction is unnecessary, as simple detection and switching will be adequate. Only if some internally redundant approach is selected, (e.g., triple-voting at subsystem or circuit level: "quadding") which puts the risk below 10^{-4} , should transfer error correction be considered. Even then it appears, the addition of correction circuits to the internal-redundancy would result in greater cost and an aesthetically less-desirable system than a dual- or multi-computer.

There is an important further consideration relative to error detection, even for single computers in the under-1000-hr MTBF range. At times other than those of peak demand, errors of commission could occur. The effects of these might range from embarrassing (generation of spurious alerts) to frightening (arming of an interceptor). Simple parity checking on transfers, plus the other "easy" checks, provide inexpensive insurance against such errors. Further, as transient error rates are unpredictable, detection permits tallying and indication of situations where transient rates are increasing, indicating impending faults. Furthermore, all error checks are valuable in the process of fault detection and location, and probably more than just the "easiest" will be included for this purpose.

Programmed error correction may be performed, but extra bits must be added to storage and registers (for the same precision), space must be added to storage for

the program itself, and, for peak-demand situations, the speed must be increased to accommodate the extra program steps.

Chapter IV

MULTIPLE COMPUTERS FOR RELIABILITY

1. INTRODUCTION

One method of obtaining reliability (for a price) is to buy extra computers. Then, quite simply, when one fails, a spare takes its place (if the spare is working). If the probability of at least one machine being available is high enough, and if there are ways of detecting failures so that repair can be initiated, the multiple computer concept becomes very attractive. The following examples use several computers where the actual problem requires but one: the extra machines provide back-up.

Several distinctions must first be made. First, there is the difference between "on-line" and "non-on-line" operation. On-line means that all the computers are operating; i.e., all the reserve machines are in the same environment as the one which is doing the work. This is unlike most "spare parts" situations where the spares are on the shelf (hence, not subject to wearout). In the following analysis, the on-line situation is assumed. This case is selected primarily because the increased reliability resulting from the off-line condition does not offset the continuous error-checking ability available in the on-line condition; and also, it is unlikely that anybody would allow a large and expensive digital computer to remain idle.

Second is the question of how service is apportioned. It might be assumed that just a fixed-service capability exists and the computers are repaired sequentially if more than one fails. On the other hand, a flexible amount of

service can be assumed wherein all the required repairs are performed in parallel; both cases will be examined here.

Consistent with the all on-line mode of operation, adequate error detection will be assumed. The case of three or more computers whose outputs are majority-voted (see Sec. A-13) will not be treated here, since the duplex (two computers) method, each with sufficient self-checking, is far superior.

As in Chap. III, only the pertinent results of Appendix A will be given. The primary problem is to ascertain the asymptotic availability of the duplex and multi-processor systems and to compare the results with the redundant and non-redundant single computers.

2. THE MULTI-PROCESSOR

The multi-processor is a way of building very large, very fast, computing systems for use in solving problems which, although very large, need not be processed sequentially. That is, certain parts of the problem can be worked on at the same time; then, perhaps, the results merged and again more processing done in parallel.

To accomplish this feat, the notion of a single, complete computer is abandoned. Instead we take a number of memory units, another collection of arithmetic units, some control units, and enough input/output machinery. If at least one of each of these units is connected together, a single computer will result. Assume now that there is enough switching equipment to connect enough units together so as to form m separate computers. Also, let there be additional switching to interconnect these m

computers so that they may communicate with each other under program control (see Fig. A-28). This is what is currently called a multi-processor. A further discussion of such machines may be found in Ref. 1.

As usual, the analysis appears in Appendix A-14, and the concern here is to compare the multi-processor with other systems.

3. SYSTEM COMPARISONS

The results derived in Chaps. II and III, and Appendix A permit some comparisons. It is impossible to compare all variations of the systems which have been discussed, but the reader can make his own comparisons for cases of particular interest.

Figures IV-1 to IV-4 compare the asymptotic availability of five data processors as the service rate is varied. These graphs are virtually self-explanatory; in all cases, the multi-processor provides the highest availability and a duplex arrangement the next highest. The exception to this statement is for very low service rates (Fig. IV-4, $\mu = .025$) where the availability is higher for duplex and non-redundant machines which are very large [$\log_{10}(N\lambda \times 10^5) > 4.4$]. As might be expected, redundant systems are never serious contenders, if service is available.

The proper choice of $N\lambda$ must again be made by the reader, but some typical systems are listed in Table IV-1. Judicious use of Tables II-5 to II-13 allow a sufficiently accurate estimate of $N\lambda$ for any proposed system.

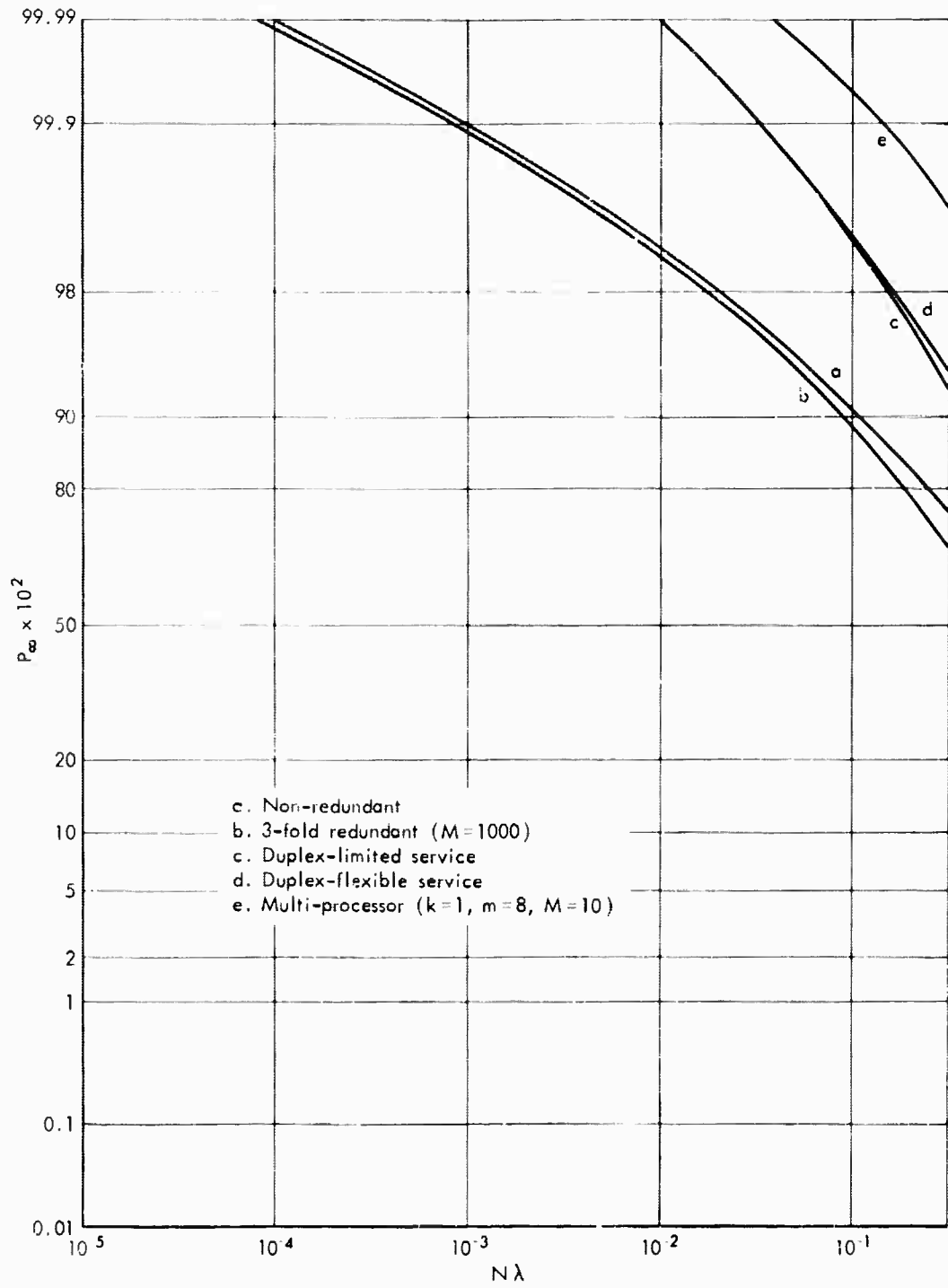


Fig.IV-1 — Comparison of systems ($\mu=1.0$)

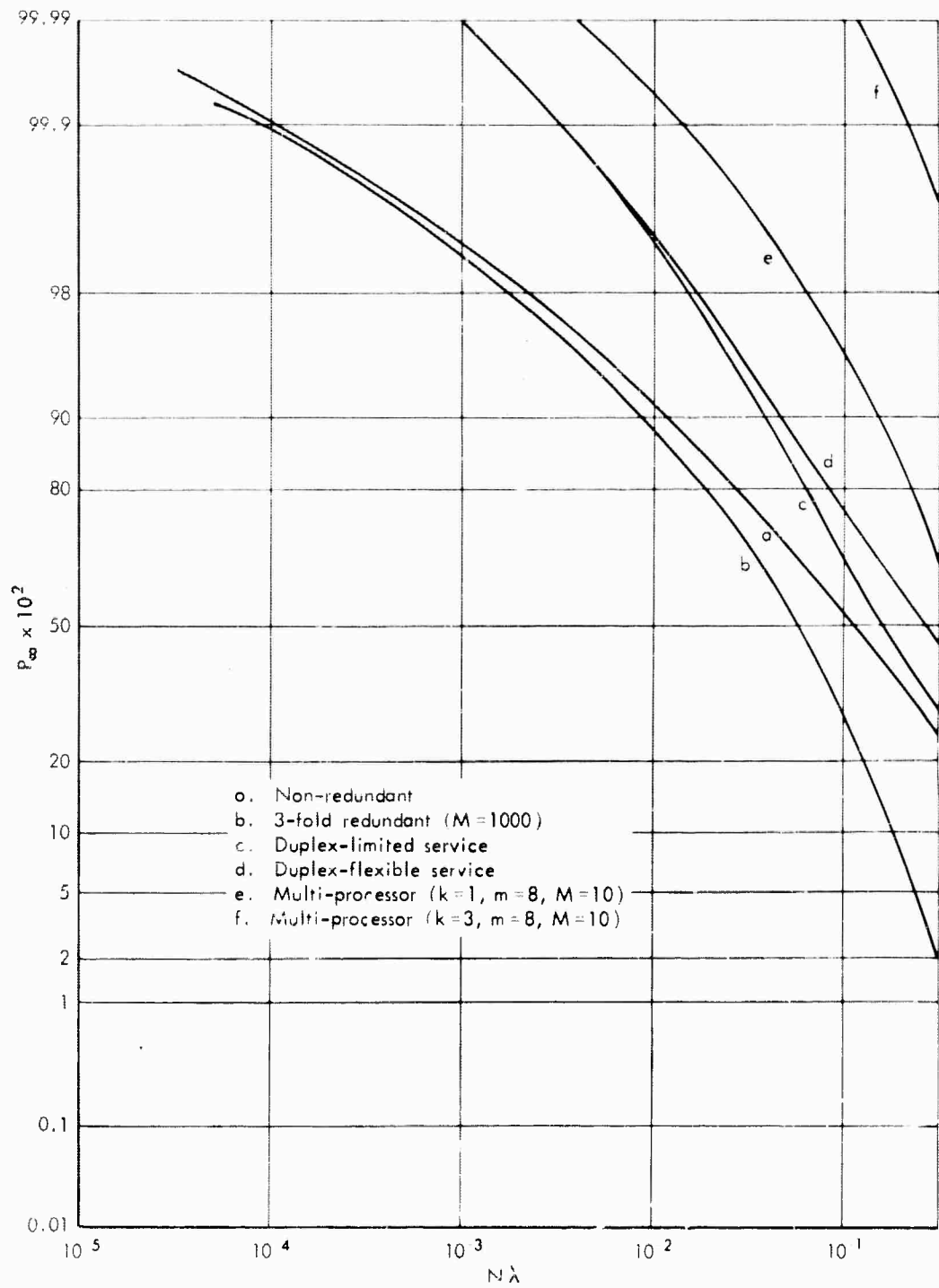


Fig.IV-2 — Comparison of systems ($\mu=0.1$)

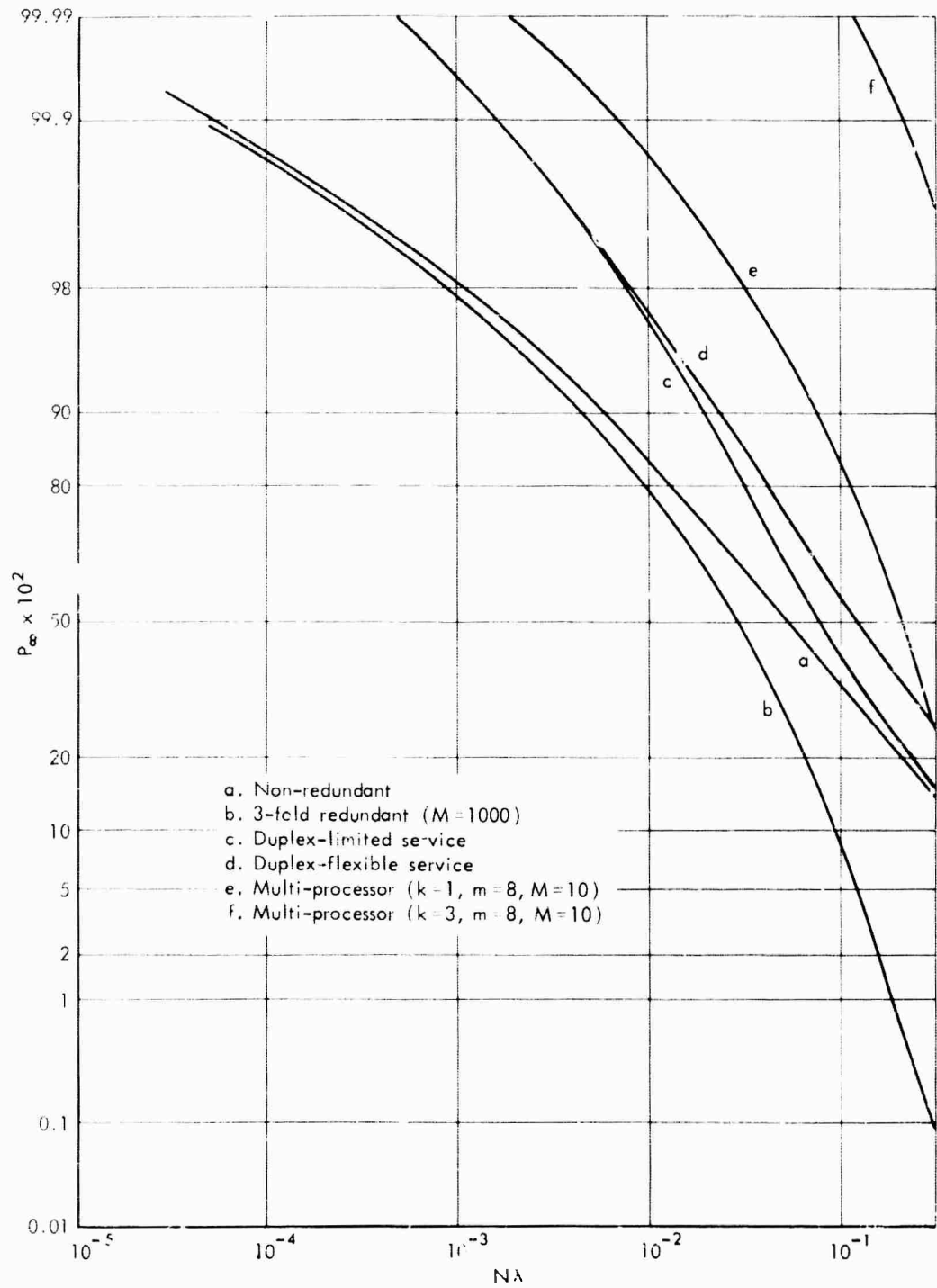


Fig. 1V-3 — Comparison of systems ($\mu=0.5$)

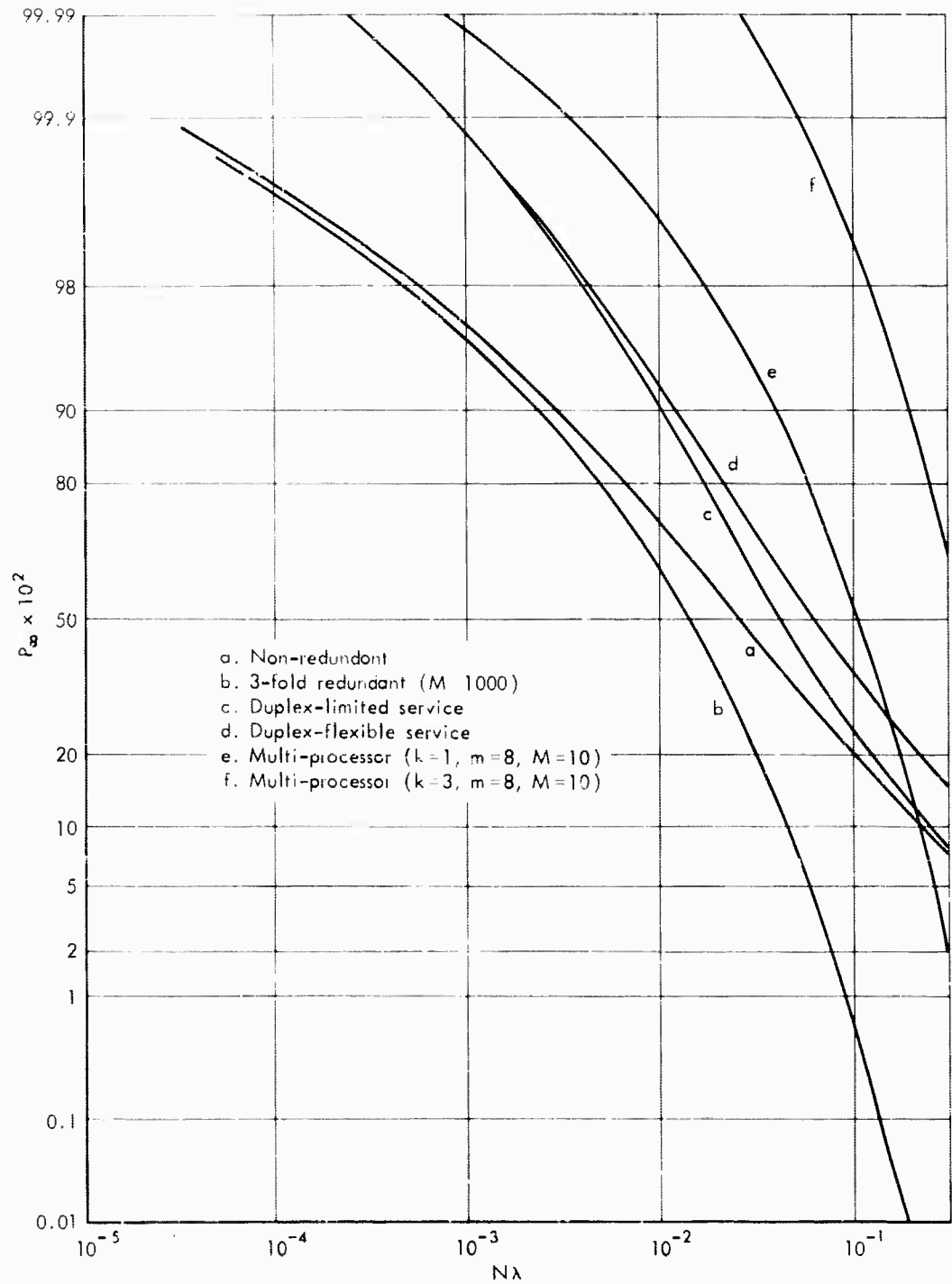


Fig.IV-4 — Comparison of systems ($\mu = .025$)

The multi-processor has its drawbacks, however. The problem of switching has been alluded to, but not treated here. The reader, when choosing a value of N for a particular multi-processor, would be well-advised to add as much as 25 per cent to his figure to allow for the added complexity in switching. We again remind the reader that the potential of the multi-processor is realized only for problems which can be partitioned into sub-problems, each capable of being worked on at the same time. Without this feature of parallelism, the beauty of the multi-processor is lost.

The problems of programming the multi-processor are also much more complicated than those of the more conventional machines. This subject will receive considerable attention in the next decade.

Table IV-1

SIZE-RELIABILITY FACTOR FOR VARIOUS PART GRADES

Parts	Size (Unit Complements)	$N\lambda$
Present computer grade	10	1.83×10^{-1}
Present computer grade	1	1.83×10^{-2}
Present computer grade with burn-in	10	$.92 \times 10^{-1}$
Present computer grade with burn-in	1	$.92 \times 10^{-2}$
Present high reliability	10	$.55 \times 10^{-1}$
Present high reliability	1	$.55 \times 10^{-2}$
Predicted 1968 computer grade with burn-in	1	1.83×10^{-4}

Chapter V

SYSTEMS CONSIDERATIONS

1. INTRODUCTION

This chapter considers a variety of topics concerning how the data processor is matched to its environment. Presumably, we now possess a very reliable computing machine, and the task at hand is to use it properly. Programming, interconnection and packaging, extreme environments, and quality control are some of the remaining problems. This list is surely incomplete--every systems man will have his own list of troublesome thorns. The maintenance process (which also includes error diagnostic techniques) is so important that a separate chapter will be devoted to it.

2. PROGRAMMING

How to write a correct program remains a difficult problem--and one for which no formula can be given. Programs are always reliable in the usual sense of the word, but large ones usually contain undetected errors--undetected, that is, until sometime later a particular set of circumstances takes the program into an untried branch which contains an error, with possibly catastrophic results.

"Catastrophic" is used here in the sense of some important, irreversible event which results from a programming error. Most errors do not result in such occurrences, but real-time control processes such as satellite guidance,

the control of a steel mill, or ballistic missile defense are situations where errors may, indeed, be catastrophic.

Unfortunately, although very large programs can be thoroughly tested, the number of tests which are required for an exhaustive checkout is astronomical and prohibitive. In view of this, the best that can be done is to set down a set of procedures which, if faithfully followed, will significantly reduce programming errors [1,2].[†]

Errors are usually divided into two classes: "coding errors" and "logical errors." Coding errors are those that might better be called clerical errors; the keypunch operator hits the wrong key, a card is out of order, etc. Everybody commits these errors, but fortunately very few of them go undetected. If the program is written in an assembly language[‡] (e.g., FAP or MAP), it will rarely assemble and execute in the face of a coding error. On the other hand, if a compiler is used (e.g., FORTRAN or ALGOL), these errors are a little easier to make.^{††} But coding errors which will still permit the program to compile and execute are still infrequent.

[†]For the reader who wishes to pursue this subject beyond what is given here, an excellent starting point is Hosier [2].

[‡]For the definition of an assembly program and compiler, see Haverty and Patrick [3].

^{††}For example, in FORTRAN, the following error went undetected for six months:

BØ = C + D

⋮

E = BØ + F

The zero was used in the last line instead of the letter O. FORTRAN defines BØ as zero and executes.

Much more troublesome and more frequent are the logical errors, primarily because they must be discovered (in cases where they are not glaringly obvious) by the programmer's own intimate knowledge of what the program is supposed to do. Logical errors result from the programmer's misreading or misunderstanding the programming design specification. For example, adding two numbers when they should have been subtracted, or taking a particular branch when $A < B$ instead of $A \leq B$. Just as coding errors are perhaps easier to commit in a compiler language, so logical errors are easier to make in the assembly language. All in all, however, fewer errors will be made if a compiler is used instead of writing directly in the assembly language.

As for the cure, very few substantive statements can be made. "Be careful!" and "Double check!" just aren't sufficient. Hopefully, the items listed below will add a little toward acquiring error-free programs.

Techniques for Avoiding Errors

Documentation. Complete specifications of what the program is supposed to do (Program Performance Specification) and how the program is to accomplish this goal (Program Design Specification) should be mandatory.

Furthermore, enough flowcharts and subsidiary documentation must be written so that bringing in a new programmer is not an impossible task. Changing programmers, or giving one programmer another's code, is a bad situation at best and should be avoided when at all possible.

Here, also, the use of a compiler makes a substantial difference; the effort required to successfully change programmers with just an assembly language can be prohibitively large.

Interface Communication. Means must be available by which the routines written by different programmers can be mated. This is usually done as the system grows, but some advanced planning could avoid a lot of the potential mismatch. Some of this is just a memo circulation procedure, but much more important might be the use of common data file compiling systems, such as the CL-1 Programming System and COMPOOL [4-6]. In addition to providing centralized control of the data file, these systems remove the artificial boundaries which are brought about by fixed-word-length machines. A programmer may now, for instance, extend the data stored in the first eight bits of a certain word to twelve bits and be assured that no trouble will ensue because that data word was already full.

Subroutine and System Debugging. Make sure that every subroutine and subsystem operates in its own right. This is done by writing enough extra code to check the subroutine without attaching it to a larger subsystem or the main program.

Use of Acceptance Specifications. Every subroutine and larger subsystem should be accepted only after having passed on an already-written acceptance specification. This specification must be written on the basis of the performance specification and not the design specification. This is a very common failing and has resulted in many errors going undetected. Again, this must be done at the subroutine level as well as higher in the system.

Simulation. Extensive use of simulation techniques can, and many times should, be used to aid in debugging and locating logical errors. Simulation can be performed at many levels--what is meant here is really a type of "micro" or "internal" simulation of other subroutines or subsystems which will enable a more thorough check of the program undergoing test. The next section will cover the topic of general input simulation (i.e., simulation of the real-world environment in which the system will find itself).

3. PROGRAMMED ERROR DETECTION

One might now ask whether or not subsidiary programs can be written which are used solely for error checking. More specifically, are there programs which a) detect hardware faults, or b) check the correctness of the operational program?

The answer to (a) is probably no; there is no strictly programmed method which will check the machine. The usual procedure is to run a large set of problems whose solutions are known, and then automatically compare the machine's solutions with the correct ones. Technicians run these programs during the morning checkout-and-preventive-maintenance period of most large data processors.

Another, more advanced way in which the program can aid in error checking is to have a collection of checking circuits which are switched to different parts of the system under program control, thus reducing the total amount of checking circuitry required. This method has been utilized in the No. 1 Electronic Switching System [7].

Much less can be done in writing programs which check other programs. By this is meant the execution of a set of instructions whose purpose it is to verify the correctness of another set of instructions. It appears to be a difficult and relatively uninvestigated area, and most opinions are not very optimistic about its future.

A simulation scheme which exercises the entire program is a possibility. Such simulations are, themselves, large programs and bring with them the problems of knowing when the reactions of the operational program are wrong. Obviously, the correct value of every variable cannot be supplied for every input, particularly when the simulation inputs are generated by some random process. "Reasonable" bounds could be provided, however, and the program scrutinized if any variable fell outside its bounds.

4. INTERCONNECTION AND PACKAGING RELIABILITY

Interconnection may be roughly classified in order of increasing complexity of the disconnection process.

A first general category includes connections intentionally designed for occasional or relatively frequent disconnection, such as for installation, removal, and maintenance purposes. Multipin plugs and receptacles and the various forms of etched card connectors fall in the group. Multipin plugs and receptacles are sufficiently familiar as to need no further description. Female connectors for etched boards (sockets) vary widely in detail design, but all attempt to combine some form of high-pressure wiping contact with easy insertion-withdrawal characteristics. Male etched board connectors (plugs) may be integral or attached. Integral, or "self," connectors are tabs of the metallic laminate, and may be redundant

(duplicated on both sides) or non-redundant, requiring like characteristics of the socket. Attached etched-board plugs may be individual male contact parts, staked and soldered or otherwise attached, or complete plug assemblies fastened to the board.

A second category of interconnections is designed to facilitate infrequent disconnection, such as for correction of assembly errors or incorporation of changes, while maintaining good long-term performance characteristics. This group includes wire-wrap and taper-pin techniques. Wire-wrap involves wrapping several turns of bare wire around a terminal, preferably of rectangular cross-section, with a special tool. The high pressure action at the corners of the terminal produces a bond in which solid state diffusion may actually take place, causing eventual improvement of the connection [8,9]. Taper-pin connections are made by first staking a tapered sleeve to the end of the wire, then driving the sleeve into a mating female sleeve with a special impact tool.

Soldering, welding, and thermocompression bonding comprise a class of techniques still permitting disconnection and reconnection, but at less convenience than the two methods above.

Finally, plating, metallic deposition, solid-state diffusion, and similar techniques can produce interconnections having specific desirable characteristics, but which are essentially non-alterable.

As part reliability increases, so does the likelihood of emergence of interconnection unreliability as the dominant contributor to failure.

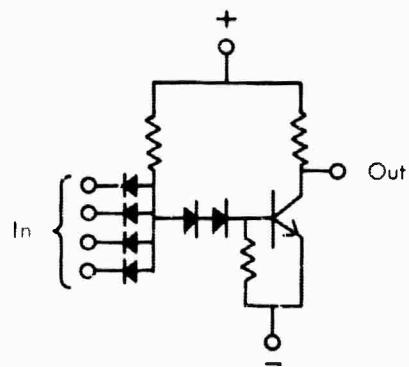
Figure V-1 shows the schematic and a possible etched board and monolithic integrated circuit layout for a four-input diode NAND circuit. The component count is six diodes, three resistors, and one transistor. Table V-1 shows the interconnection count, assuming planar transistors and diodes in individual cases, and a conventional can closure for the microcircuit.

Table V-1

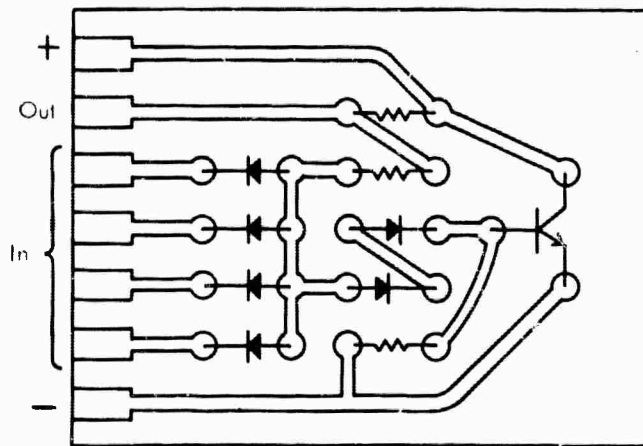
CONNECTIONS FOR ETCHED BOARD AND INTEGRATED CIRCUITS

Type	Etched (Discrete)	Integrated
Die attach, bond, weld	23	15
Diffused aluminum contact	8	24
Solder joint	21	0
Diffused aluminum path	0	10
Copper path	11	0
External (socket, solder, weld)	7	7
Total, all types	70	56

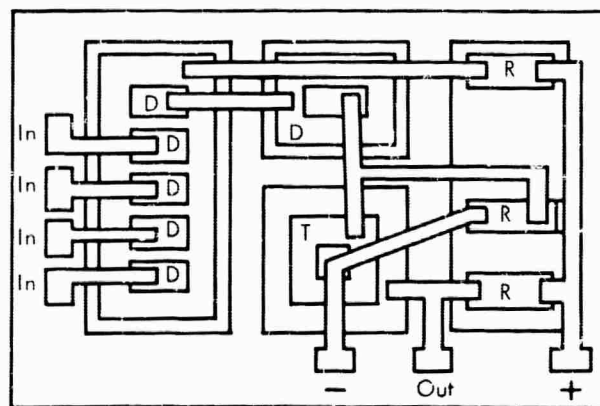
For discrete circuits, the parts reliability estimates cover the first two items (31 instances). For integrated circuits, all but the last item are similarly covered. Failure rates may be estimated as .020 for the discrete circuit (four equivalent transistors @ .005) and .010 for the integrated circuit. If interconnection reliability, to be neglected, must be an order of magnitude better than



a. Schematic



b. Etched board layout (2x to 4x actual size)



c. Integrated circuit layout (actual size --> □)

D = diode
T = transistor
R = resistor

Fig.V-1 —Etched board and monolithic integrated circuit layout

that of the parts, then individual connection reliabilities must be (ignoring the copper paths):

$.002/28 = .00007$ for the discrete circuit;

$.001/7 = .00014$ for the integrated circuit.

With suitable care these figures are apparently achievable for wire wrap, soldered, and welded interconnections, as indicated in Appendix E. Taper pin connections are probably somewhat less reliable, but are becoming unpopular due to topological limitations and the impossibility of automatic assembly. Connector data obtained from Ref. 10 is probably misleading for modern etched-board connectors. One source [7] indicates the feasibility of etched-board connector failure rates completely compatible with best semiconductor device rates, using only non-redundant self-connectors. Although brute-force life-test verification is not, and may never be, available, one may reasonably assume, on the available evidence, that deposited aluminum interconnections and thermocompression bonds (again with the proper control) can achieve reliabilities equal to or better than the more gross methods. One technique under development uses all-diffused interconnections, mounting silicon integrated circuits on a silicon "carrier board." Potential reliability is clearly as good as that of the processes forming the active elements themselves, though reduction to practice may present considerable difficulty.

The degradation or failure indication of an interconnection is, of course, increased resistance, which in

turn results from increased resistivity or decreased cross-section of the conductive path. Electrochemical and thermochemical changes such as electrolysis, oxidation, and corrosion may either increase resistivity or decrease cross-section. Partial detachment through physical stress decreases cross-section. The solder joint requires more care than the others in that cleaning is required to remove flux residues.

It should be noted that failure ascribed to interconnections may result either from static or dynamic effects--it may be caused by the interconnection design, as well as interconnection degradation. The problems of interconnection design for optimum dynamic behavior have been covered in the literature [11], but, with respect to noise considerations, each new system requires its own analysis.

5. EXTREME ENVIRONMENTS

Packaging reliability, for ground-based equipment in a controlled atmosphere, need not be given detailed consideration here, unless protection against physical shock from nearby explosions is required. The parameters of the shock may be estimated, and conventional shock-isolation procedures incorporated. Protection against short-term high-intensity radiation involves consideration of two effects--permanent damage, and transient behavior [12-14].

With respect to permanent damage, the type, intensity, and duration of radiation becomes a degradation/failure stress factor to be incorporated in part selection, assignment of end-of-life design parameter limits, and estimation

of failure rate. If the estimated failure rate or spread of degradation limits for available part-types is intolerable, some form of shielding will be required.

A transient error or fault resulting from radiation is no different from that caused by any other unanticipated stress; all previous comments on this subject apply.

6. THE ROLE OF THE MANUFACTURER

A final system consideration is essentially the summation and reiteration of a number of allusions elsewhere in this report--specifically, the integrity of the manufacturer. An entire organization might possess an attitude which operates essentially to the detriment of the product. Conversely, organizations exist which evidence a technical and administrative orientation which tends to guarantee production of outstanding systems. No method of computing mean-time-between-failure can take into account intangibles such as

- o Lack of firm technical policies in design, materials selection, handling, and assembly, and checkout;
- o Lack of understanding of, or agreement with, said policies on the part of the technical team;
- o Emphasis on quantity, rather than quality, in engineering personnel--often the result of an over-aggressive expansion-oriented, marketing policy;
- o Misunderstanding of the roles of inspection, quality control, production control, reliability, and PERT, and compounding of the felony by the introduction of even more esoteric groups and philosophies (often called "testing quality into the product").

Chapter VI

THE PROCESS OF MAINTENANCE

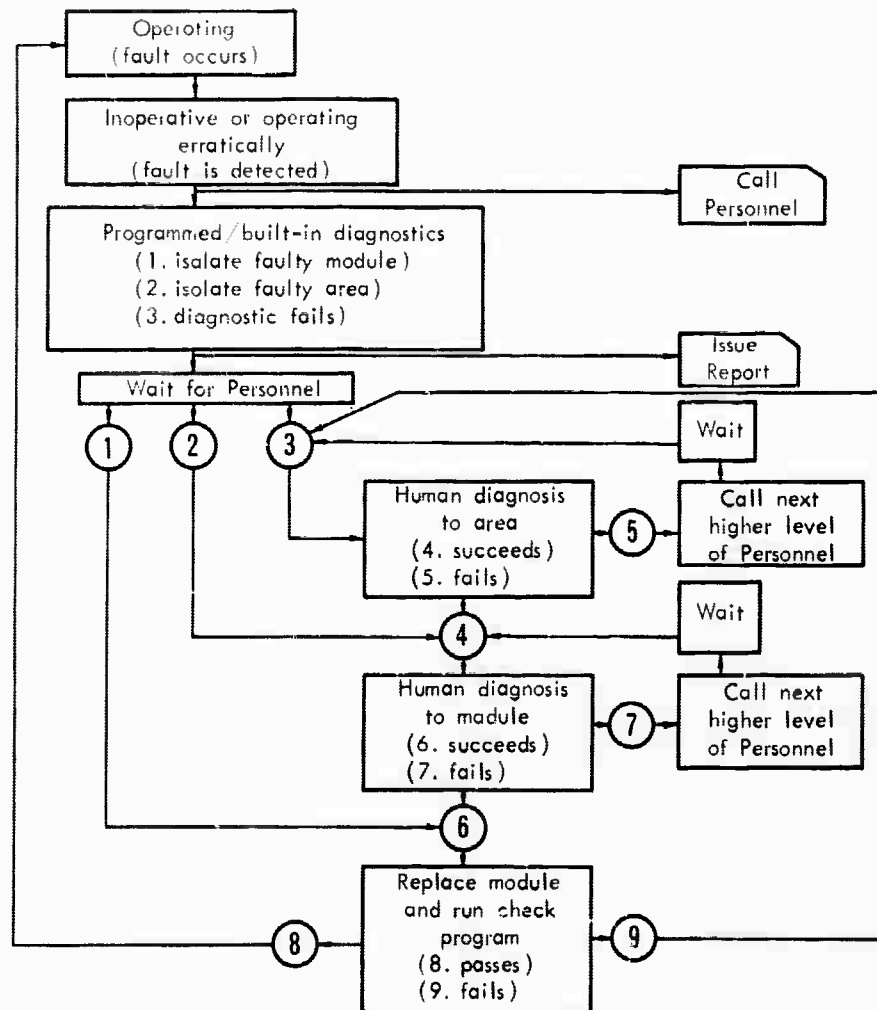
1. INTRODUCTION

Whether downtime starts with a failure or the detection of failure is a philosophic distinction similar to that outlined at the start of Chap. II. Ideally, detection of failure would always precede any erratic output action, and approaches to this ideal have been discussed. This chapter essentially concerns the period between detection of failure and restoration of the system to on-line operational capability--the "corrective maintenance period." Corrective maintenance differs from preventive maintenance in that it is unscheduled. It occurs whenever the fault is detected rather than when it (maintenance) is scheduled.

Figure VI-1 is a flowchart of the field corrective maintenance process. The time from (FAULT IS DETECTED) to re-entry into the operating state is called the "time to repair" (TTR), or, better though rare, the "time to restore." The time from (FAULT OCCURS) to (FAULT IS DETECTED) might be called the "time to detect," and will be assumed negligibly small relative to any TTR or time between failures (TBF).

The chart contains some implicit simplifying assumptions: a) The diagnostic period always ends prior to arrival of personnel, and b) human diagnosis always follows the orderly pattern of system-to-group-to-module.

The various actions in the chart will be considered sequentially, commencing with fault detection. Detection of a permanent error by any means, as described in Sec. V-3



(Expressions in parentheses are exit conditions)

Fig.VI-1—Flow chart of field corrective maintenance process

is, of course, equivalent to detection of a fault. Special built-in fault detection circuits are available which permit more rapid detection of situations which, if unchecked, may lead to multiple failures (such as overvoltage-undervoltage relays on power supplies and temperature interlocks).

Direct programmed fault detection is possible only if some built-in provisions are made to convert electrical signals to a form accessible as computer information. The usual problem exists in the case of a control fault which does not permit the program to run.

Diagnostic (fault location) methods may also use a combination of built-in and programmed implementation. A built-in fault detection scheme often provides location automatically.

2. DIAGNOSTIC TECHNIQUES

Success in detecting and locating faults at an operational site depends on a number of factors:

- o The time available to perform the necessary functions;
- o The built-in features of the system which aid in the diagnosis;
- o The level to which fault location is required;
- o The quality of the personnel.

Totally automatic approaches that will detect and locate any possible system fault down to the desired replaceable/repairable level are not yet available.

However, a number of automatic approaches are in use that will effectively handle a major sub-set of the possible failures.

A conventionally designed computer requires a large portion of the computer to be in working condition if the computer is to be capable of executing a program. This portion is usually 30 per cent or more of the computer circuits.[†]

Therefore, when developing a computer program that will automatically diagnose failure (at least in the computer main frame), one must remember that a large portion of the machine must be in working order merely to run the program correctly. Moreover, if a fault exists in the necessary basic circuits, nothing more than a very gross isolation of the failure is possible.

Another drawback of conventional computers is the inaccessibility of the various registers. In general, the only registers available for comparison purposes are the accumulator and multiplier/quotient registers. The contents of all other registers must be deduced. Of course, most machine consoles have display capabilities for these other registers, but are not usually accessible to the machine's comparative circuits. Furthermore, one cannot set the internal state of a conventional computer from an external location (e.g., the console). Generally, the ability to set the internal state is limited to that which is required

[†]The value 30 per cent applies to conventional designs. As shown later, significant decreases in this fraction (usually called the "hardcore") can be achieved.

for program running and monitoring, and the hardware is assumed to be functioning properly. Hence, the state of the machine after any operation is unknown unless the machine is operating without errors.

Three basic approaches to diagnostic programs for computer mainframes exist. Generally speaking, the part of the program devoted to fault detection is very similar in all three cases. The methods differ primarily in the techniques used for fault location and the extent to which this is achieved automatically.

Method 1

The first approach is the one most commonly used in both commercial and military installations. This type of program is characterized by the following features and assumptions:

- o The set of circuits functioning properly is sufficient to enable the program to operate;
- o Only single faults will occur;
- o Fault isolation can be deduced from a set of fault indications for all possible, acceptable failures (i.e., it rarely says what the indications of the failure were);
- o The technical level of the technicians and the time available is adequate to deduce failure locations;
- o Few or no additional circuits are included in the system to perform fault location;
- o It does not attempt to exhaustively test the logic; rather, it performs a "sufficient" number of tests with "randomly" generated data and specially designed sample patterns.

A description of a program of this type is given by Bashow [1].

This approach is probably inadequate for field-operational data processors, as it requires extremely competent and well-trained technicians for effective application.

Moreover, about 30 per cent of the computer is not treated at all by this approach. Since the circuits in that 30 per cent are not basically different from the rest of the system, about 30 per cent of the faults can not be treated without the addition of extensive testing equipment.

This approach is also characterized by the lack of knowledge of the effectiveness of the procedure. The test routines are heuristically developed and no methods exist to test their completeness. This would lead to continuous debugging and updating of the program even in the absence of the inevitable system modifications.

The results claimed by Bashow [1] clearly underscore the inadequacy of this application for field applications. Of 34 actual tests, only 23 were adequately diagnosed by the program. No further statistics were given as to the efficacy of the program.

Method 2

The second approach, used to some extent, utilizes a second computer to test the first computer. The following features characterize it.

- o Fault detection and isolation is totally automatic (to whatever extent the program is capable of such detection and isolation).

- o A great deal of pre-processing is required to accumulate the data necessary to realize detailed fault isolation. This is accomplished both by simulation and by actual introduction of failures.
- o Some additional hardware is required in the computer to effect the interconnection of the two computers.
- o Generally, only single faults are assumed to exist, although multiple faults could be handled.
- o Each fault is associated with a given, unique sequence of failure indications that can be monitored by the "good" computer.
- o It attempts to exhaustively test the logic and to, at least, detect all possible errors.
- o It assumes that the technical level of the technicians and/or the time available is low.

The multiple computer approach is more suited to field diagnosis. However, it presents several problems. First is the cost of generating the fault-location data and of maintaining the data (i.e., keeping it current with state of the system).

Another problem is determining which faults to consider. Once the set of acceptable faults is selected, the entire structure of the maintenance procedure is fixed. No other faults can be treated, and any other fault may be improperly diagnosed as one belonging to the acceptable set. This could lead to great difficulties in actual maintenance because, although the average repair time might be very low, the maximum repair time is likely to be extremely high, perhaps a matter of days.

Tsiang and Ulrich [2] give the following statistics: For a computer containing 8000 circuit packages having 6500 transistors and 45,500 diodes, the diagnostic program required 7200 words, and 50,000 faults were treated. Only single failures were considered.

For 75 per cent of the included faults, location is to one circuit package and for 13 per cent, location is to two packages. Probably 75 per cent of the failures are locatable with this method.

The fault diagnostic information was generated by actually creating the failure in the hardware rather than by simulation. The generated data consisted of about 60 million bits and was reduced to a dictionary of 1290 11x15 pages. The project required about 13 man-years and 250 hr of machine time.

Method 3

This rare approach is characterized by the following features:

- o The computer diagnoses itself; however, additional circuits are provided to set and reset all (or most) storage devices "directly," to read the contents of all (or most) storage devices "directly," and to provide alternate control paths;
- o A fairly small percentage of the circuits (~10 per cent) must be functioning properly to allow the program to operate;
- o Fault detection is totally automatic, but fault location is only partially automatic;
- o A great amount of pre-processing is required to accumulate data for automatic fault location;

- o The approach can handle all single faults and some multiple faults;
- o Faults are located by observing the point at which the test failed (gross isolation) and the resulting indications, and then deducing the cause by referring to appropriate documents;
- o Testing and, consequently, fault detection are exhaustive;
- o It assumes a fair level of technical competence from the technicians and a fair amount of available time.

This method has many features of merit. A description of a system of this type is given by Carter [3]. The powerful control capability that can be exercised over the computer is particularly attractive. It allows the development of simple test procedures by sub-dividing the system into manageable chunks. Most of the tests and the data needed to analyze the results can be automatically generated from the data stored in design automation files. This leads to a minimum of transcription errors, and completeness in the testing procedures. It also allows automatic updating of diagnostics in response to system changes.

The method has drawbacks, however. The hard core which cannot be reached by the program has been reduced to ten per cent, but this ten per cent still must be handled by other means. The enormous cost must also be considered a drawback, if not an absolute deterrent. Over 100,000 instructions are required for the IBM System/360 and the cost might conservatively be \$500,000. Finally, fairly sophisticated technicians will be needed to implement this method.

Further consideration here of these methods is unprofitable. The nature of the service, the allowable budget, and the value of very low repair times all must be estimated before recommendations can be made.

Much work has been done on error diagnosis and fault location which cannot be reported here. The general theory of diagnosis in switching and sequential circuits is discussed in Refs. 4-8. The subject of transient and intermittent errors and the ability to re-examine the program has not been considered, but information may be found in Ref. 3. Elegant methods such as those being developed for IBM System/360 have a possible dividend in that logical design errors are also sometimes revealed [3,4].

3. THE MAINTENANCE MODULE

The process of fault location is significantly affected by the size of the maintenance module. If the module is an entire computer, fault location becomes identical with detection. This creates a rather bulky and expensive module, however. At the other extreme is the single part as a maintenance module. Fault location to a single part at field level, and the problems associated with installation, offset the portability and economy of the single-part module.

As an attempt to evaluate the tradeoffs involved in module size selection, consider a computer composed of

- 100,000 transistors of 25 types,
- 150,000 diodes of 25 types,
- 150,000 composition resistors of 50 types,
- 50,000 film resistors of 50 types,
- 50,000 mica capacitors of 50 types.

The total is 500,000 parts of 200 types. An additional 200 types might be required in other categories (paper and electrolytic capacitors, zener diodes, transformers), making 400 types altogether.

When parts are combined with plug-in modules, they will be partitioned into various numbers of modules of various types, which will determine the nature of a complete set of spares. Although module partitioning statistics are not readily available, the few instances discovered led to formation of the following highly hypothetical partitioning rules:

- o Half of the modules of the system are replaceable by 16 types;
- o Half of the remaining modules are replaceable by 16 more types;
- o And so on, until the remainder is less than 32, when all remaining modules are assumed to be unique.

For the 500,000-part computer, some partitioning options are listed in Table VI-1, assuming an average part cost of \$1.00. For the one-part module, the partitioning rule gives 238 types, but the assumed number of part types was used instead. Note that "cost" includes parts cost only. The more complex modules would be more expensive, due to assembly labor and burden.

The most significant factor to balance against cost of a set of spare parts is time to diagnose to the maintenance module level. Estimates of diagnostic time vs. maintenance module size are at least as nebulous as the partitioning rule, but one possibility is to hypothesize

Table VI-1

SPARE PARTS COST VS. NUMBER OF MODULES

Modules	Module Types	Parts/Module (=Cost in \$)	Parts Cost of Spare Set, \$
1	1	500,000	500,000
5	5	100,000	500,000
50	41	10,000	410,000
500	79	1,000	79,000
5,000	102	100	10,200
50,000	185	10	1,850
500,000	400	1	400

that diagnostic times are related as the logarithm (base 10) of the module size ratio. This yields the estimate shown in Table VI-2.

Figure VI-2 is a plot of the relative cost of a set of spares, and relative time to diagnose. Also shown is an equal-weighted sum, indicating that for any partitioning rule, diagnosis time rule, or weighting of cost-of-spares versus cost-of-diagnostic-time, there should be some optimum region of maintenance module size.

A final factor in the corrective maintenance flow is the nature and amount of test equipment for field diagnosis, and the extent to which system, logic, and circuit design facilitate this activity. In early computers, this system aspect was entirely neglected on the assumption that a clever field engineer with an oscilloscope could, sooner or later, figure out what was wrong.

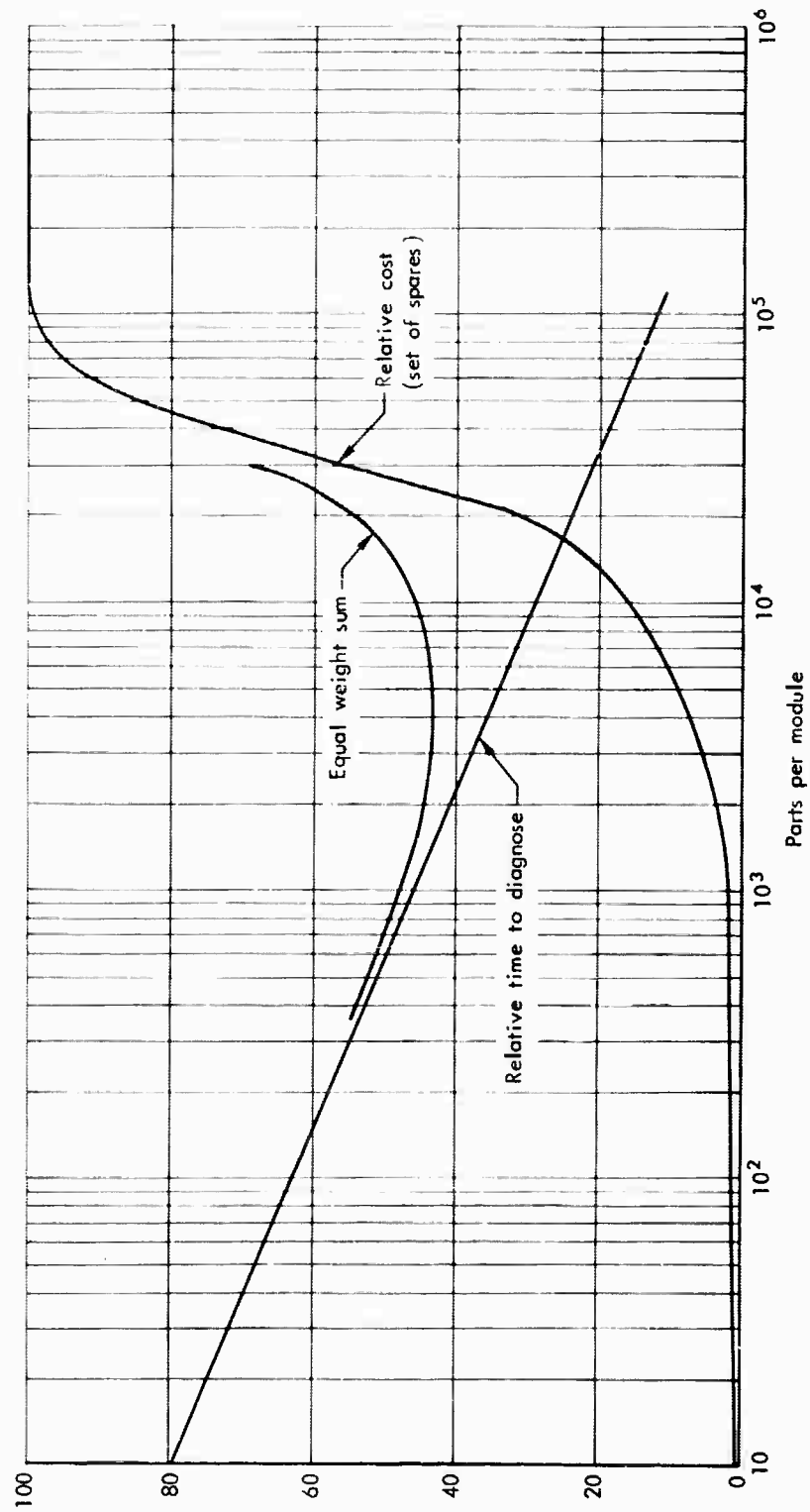


Fig. VI-2—Relative cost and diagnosis time versus parts per module

Table VI-2
ESTIMATED DIAGNOSTIC TIME VS. NUMBER OF MODULES

Parts in Module	Relative Diagnostic Time	Per cent of Maximum
1	5.7	100
10	4.7	82.5
100	3.7	65.0
1,000	2.7	97.4
10,000	1.7	29.8
100,000	0.7	12.3
500,000	0.0	0.0

Currently, all enlightened digital systems manufacturers admit the need for compromising system design relative to operational specifications, to permit rapid, hazard-free field diagnosis. Requirements extend from provision for connection of test equipment to make "accidents" impossible, to design of specialized analysis equipment brought into the site when self-diagnosis fails. As this is usually the case when control faults occur, system provisions for forcing control states and sequences as well as for setting register contents are usually included. Many other features of various degrees of sophistication may be provided, assuming some sort of cost/risk tradeoff evaluation underlying each decision.

The topics of depot maintenance, spare parts logistics, and personnel training are significant, but are beyond the scope of this report.

4. PREVENTIVE MAINTENANCE

In most practical systems, certain items of equipment are necessary, yet have failure rates so high that they completely dominate the system failure behavior. Notable examples are magnetic tape transports and typewriters. For systems requiring extreme reliability, three solutions to this problem are available:

- o Eliminate such devices from the on-line portion of the system, relying on them only for peripheral functions not essential to the primary objectives;
- o Provide some form of redundancy;
- o Design a preventive maintenance plan.

Preventive maintenance introduces scheduled periods of downtime of controlled duration to decrease instances of non-scheduled periods of downtime, with a net saving in total downtime. Preventive maintenance periods initiate certain standard procedures such as adjustments, cleaning, and lubrication, as well as specially designed tests which indicate the relative degradation of some parts. This latter process, called marginal checking, permits replacement of parts showing anomalous degradation before failure actually occurs. If the marginal check causes a "weak" part to fail, it is really a form of latter-day burn-in. More often, the effect of a marginal check is to introduce a transient fault which is then diagnosed in the usual

manner. Typical marginal checks for solid-state digital systems include:

- o Increase of clock rate;
- o Decrease of a voltage to which the noise immunity is monotonically related;
- o Decrease of a voltage to which the logic signal swing is monotonically related;
- o Introduction of simulated transistor I_{CBO} ;
- o Simulated reduction of transistor h_{FE} .

The effectiveness of any such measures depends on the expected degradation behavior and the response of the parts to the simulated enhancement. If occasional voltage breakdown were a failure mode, raising the supply voltage as a marginal check would be a disastrous election, due to the severely nonlinear behavior near breakdown.

Clock rate increase is easily accomplished, either continuously or in one or more increments. The second and third items depend on the exact design--that is, the availability of a single voltage supply exhibiting the required relationship. In some systems, the third and fourth items above are accomplished by adding a resistor to the base of each transistor. The free ends of the resistors may be connected together by a single marginal check bus which, for one polarity of applied voltage injects artificial I_{CBO} , and for the other polarity "robs" drive current, thus requiring higher h_{FE} for performance. The penalties for this approach are a) added stray capacitance at the base for high-speed circuits, and b) loss of current and possible intercircuit coupling when in the normal operating mode.

Another approach is to remove maintenance modules from the computer and test them individually on a marginal checker. This practice should be strongly discouraged, as tampering in this way with an otherwise operative system is probably more conducive to failure than normal part behavior.

The "cleaning, adjustment, and lubrication" type of preventive maintenance, where required, must obviously be provided. Marginal checking should be provided only if justified by overall maintenance logistics, including the effects of multiplicity of systems or subsystems for redundancy.

All the above approaches (except module removal), may be carried out automatically, permitting fast return to the operational mode, on demand.

VII. SUMMARY

1. INTRODUCTION

This section summarizes briefly but quantitatively the contents of this report, while the Introduction (Chap. I) presents the motivation for the work, gives the scope of the report, and defines certain key terms. The reader who requires only the conclusions and a few numbers and is willing to accept them without proof should find this summary sufficient for his needs. Topics of lesser importance, or those which are too difficult to summarize succinctly will only be cited here, and a reference to their location in this Memorandum will be given.

Tables and figures given here sometimes duplicate material which appears in the text. Where this threatens to be excessive, the reader is referred to the text for a specific figure. Recommendations are not duplicated here; again, the reader will be referred to a specific location in the report, and in all cases this summary supplies sufficient context for their understanding.

In this report, the authors attempt to estimate the availability of a data processor by starting with the smallest part and carefully investigating its failure behavior.[†] They concurrently construct a mathematical model of system availability which gives the desired results for a wide variety of systems if the failure behavior

[†]Definitions, e.g., "availability," are not repeated here; they appear in Chap. I.

of the part, the service features, and the size of the system are known. These two efforts merge and estimates of current and predicted systems availability appear as a function of the size of the processor and the type of service provided.

In addition to this central theme, developed in Chaps. II through IV and Appendices A through E, many more qualitative topics receive study in varying degrees of thoroughness. Among the more important ones are circuit design, logical design, programming, and maintenance (Chaps. III, V, VI). The purpose here is to relate these topics to machine availability and to recommend improvements.

The authors know no way of quantitatively relating circuit design, logical design, programming, or maintenance to availability. What is known, without equivocation, is that lack of care in any of these areas causes the processor to be down for unwarranted lengths of time--long after the entire system has been accepted and declared operational.

2. PARTS AND A DEFINITION OF FAILURE

After a brief description of which parts are treated (transistors, resistors, capacitors, integrated circuits, but not mechanical devices of any kind, except connectors), a discussion of failure and its many meanings is given. Paraphrasing Chap. II, we conclude that a part has failed when, under some combination of normally applied stresses, one or more parameters of the part vary in such a way that the functional group containing the part does not perform its role. The point which needs emphasis is that the concept of failure has no meaning apart from a concept of proper circuit function.

3. PART FAILURE MECHANISMS AND DISTRIBUTIONS

Chapter II continues by discussing the manner in which part parameters relate to stresses, and this leads to a discussion of the two current methods of handling reliability problems: a) statistical analysis, and b) the physics-of-failure approach.

The origin, modes, and mechanisms of failure are investigated and some conclusions are reached about the existence and nature of failure rates. These conclusions may be summarized as follows:

- o Some form of decreasing failure rate for the total part population will be observed which is in no way correlated to the predicted behavior of the ideal part.
- o The total part population shows a decreasing failure rate because, and only because, various controllably small subgroups show initially increasing failure rates of various forms until every member of the subgroup has failed, at which time the failure rate of the entire remaining population effectively decreases.
- o Eventually, a universal "wearout" mechanism (e.g., diffusion processes for solid state devices) will cause the failure of all parts. This process has an increasing failure rate, but systems normally operate so far out on the "left tail" of this distribution that it is not a factor.

Arguments are presented for the case of continuing to use the exponential failure distribution in the face of this evidence, and Chap. II goes on to estimate the parameters of the decreasing failure rate Weibull distribution from the available evidence. Those estimates give marginal confidence, to say the least.

4. STANDBY CONDITIONS

It is further concluded that there are no compelling reasons to maintain the computer (assuming it is not needed) in a quiescent, power-off condition. Certainly most of the parts would benefit from special no-power storage, but the lack of error checking and the danger from power transients suggest an advantage in keeping the system on. Last but not least, the likelihood of permitting a large, very expensive digital computer to remain inoperative in order to gain a little reliability is quite small.

5. THE "EQUIVALENT TRANSISTOR" COMPUTER AND SOME TYPICAL SYSTEMS

Assuming that a computer consists primarily of transistors, resistors, and capacitors, each with known reliability, it is possible to treat the same machine (from the standpoint of reliability) as one which is constructed entirely of transistors, by computing how many resistors (or capacitors) it takes to provide a failure rate just equal to one transistor, then exchanging parts in this ratio.

To be more accurate, more than one category of resistor and capacitor should be employed, since the failure rate of a part also depends on its use in the circuit. This complicates matters to the extent that the relation, say, between resistor failure rate and the stresses applied to the resistor in its circuit must be known in considerable detail--a nontrivial task in most cases. When this is done, the size of a computer is measured in "equivalent transistors." The complexity of

some existing systems using this measure is repeated in Table VII-1.

Table VII-1
COMPLEXITY OF EXISTING SYSTEMS

System	Complexity in Equivalent Transistors
FSQ-32	383×10^3
FSQ-31V	274×10^3
CDC-3600	97×10^3
CDC-1604A	82×10^3
Univac 1107	78×10^3
Burroughs B-5000	67×10^3
Honeywell H-1800	49×10^3
Honeywell D-825	41×10^3
SDS 9300	35×10^3
USQ-20	30×10^3
IBM 7090/44	26×10^3
GE 215/225/235	22×10^3

6. THE UNIT PARTS COMPLEMENT

To get an even more tractable way of describing the size of a computer, the concept of a unit parts complement is introduced in Chap. II. This is taken to be 10,000 transistors, 15,000 diodes, 15,000 composition resistors, 5000 film resistors, and 5000 mica capacitors. These reduce to 18,334 equivalent transistors. The size of a system is subsequently expressed in multiples of this unit parts complement.

7. PREDICTED PART FAILURE RATES

Part cost versus reliability and procurement policy is investigated next, for integrated circuits as well as the discrete type. When all is said and done, the most significant results are shown again in Table VII-2. When utilizing the availability graphs presented below, the reader should select part failure rates from this table.

8. INTEGRATED CIRCUITS

Further discussion of the reliability of integrated circuits is given, and Chap. II ends with conclusions and recommendations on part procurement and handling which really cannot be adequately summarized here.

9. CIRCUIT DESIGN

Chapter III introduces the subject of circuit design and describes the philosophy of "bogey," "worst case," and "statistical" design procedures. The latter two methods

will insure the satisfactory operation of the circuit even though the parameters of some parts vary over wide ranges.

Recommendations for good circuit design are given next; these are independent of earlier parts of the chapter and can be read by themselves. Following them is a discussion of logical design which contains some suggestions on how to minimize errors in the design process.

Table VII-2

PREDICTED PART FAILURE RATES
(%/1000 hr, 10-year average)

Part	Good 1965	Best 1965	Best 1968
Resistor (composition, metal film, tin oxide)	.0003	.0001	.00003
Capacitor (glass, mica)	.0003	.0001	.00003
Diode (silicon planar)	.005	.0015	.0005
Transistor (silicon planar)	.01	.003	.001
Integrated circuit (silicon planar)			
10 equivalent parts	.02	.005	.0015
30 equivalent parts	.04	.009	.0025
100 equivalent parts	.07	.015	.0040

10. PART REDUNDANCY

Chapter III proceeds to the concept of part redundancy to counter outright part failure. The question of whether to adopt the redundant circuit technique as a means to system reliability does not appear until Chap. IV, but Chap. III compares non-redundant with various types of redundant computers for the asymptotic case (i.e., for very large times after first turning the system on) in Fig. VII-1, and for the transient phase in Fig. VII-2. It is later shown that if service is available, circuit (or part) redundancy is not the best way to obtain high system availability. But if no service is available, redundancy is best.

11. FAILURE DETECTION

Finally, Chap. III takes up the problem of failure detection. In this chapter, failure detection and correction are considered from a circuit standpoint; later this same topic is treated from a programming viewpoint. Generally, it is concluded, error detection can and should be performed wherever possible, and the expense is not exorbitant. But error correction, while often feasible, probably does not warrant the expense, particularly in light of the elegant fault-diagnosis techniques which are currently under development.

12. MULTIPLE COMPUTERS AND THE MULTI-PROCESSOR

After defining various methods of using more than one computer to achieve higher availability (duplex, triplex, and multi-processor), these systems--the redundant and

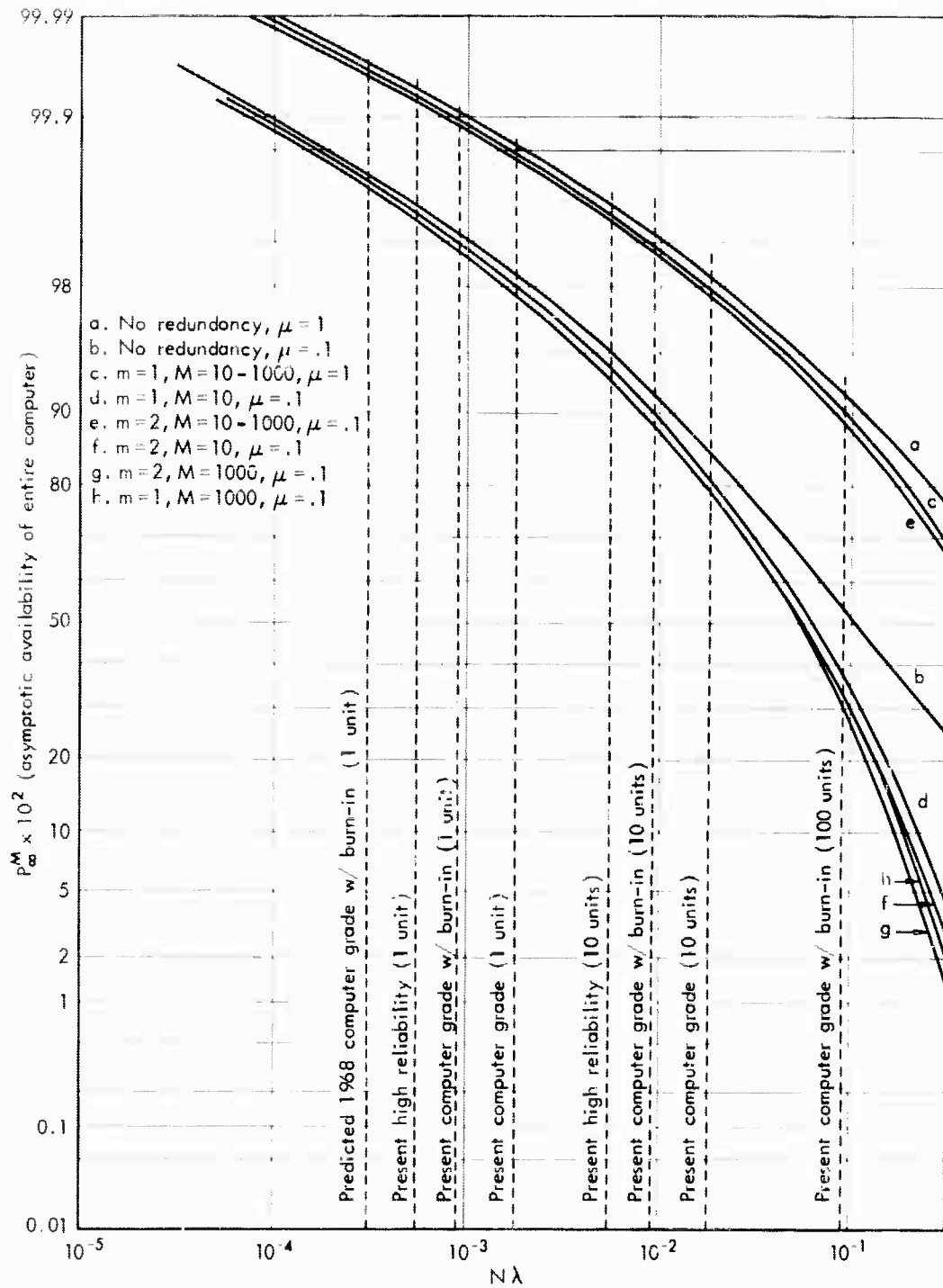


Fig. VII-1--Asymptotic availability of redundant computers (exponential service)

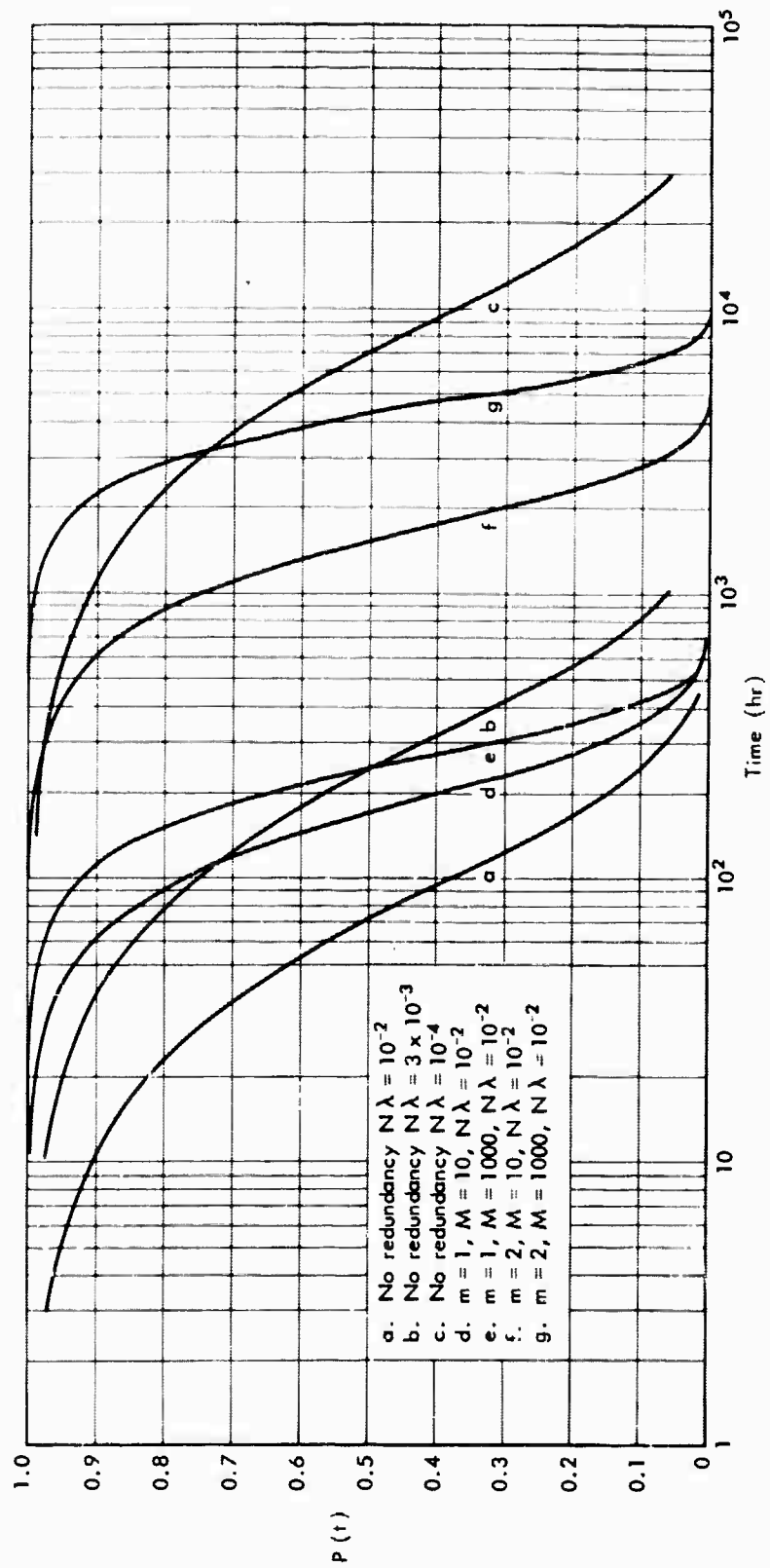


Fig. VII-2—Availability of redundant computers
(no service)

single non-redundant computers--are compared in Chap. IV. Figures VII-3 to VII-6 show the results. The list of symbols should make them self-explanatory. In almost all cases, the multi-processor is better than the duplex arrangement, and the duplex system is superior to the redundant system.[†] This is true only for the asymptotic case, which is the proper region for our attention when service can be provided. The transient behavior of multiple computers is not calculated for most cases, because this implies a well-defined, relatively short, system lifetime with no service (which just is not the case with most ground-based systems). Unquestionably, non-redundant methods will not measure up to redundant ones during this transient phase.

Furthermore, the problem of the increased size and complexity of the program must be carefully investigated before making a decision in favor of the multi-processor. This decision is so problem-dependent, and the entire subject so new, that useful guide lines are at this time impossible.

13. PROGRAMMING

Chapter V points out that methods cannot be presented that insure perfect programs. After all, correct codes are

[†] Recommending the multi-processor requires restraint, because the excellence of this system appears only when it solves problems which allow simultaneous processing of different parts of the problem. Problems which, although very large, must be done in a strictly sequential manner are not candidates for the multi-processor, and this system is no longer first choice.

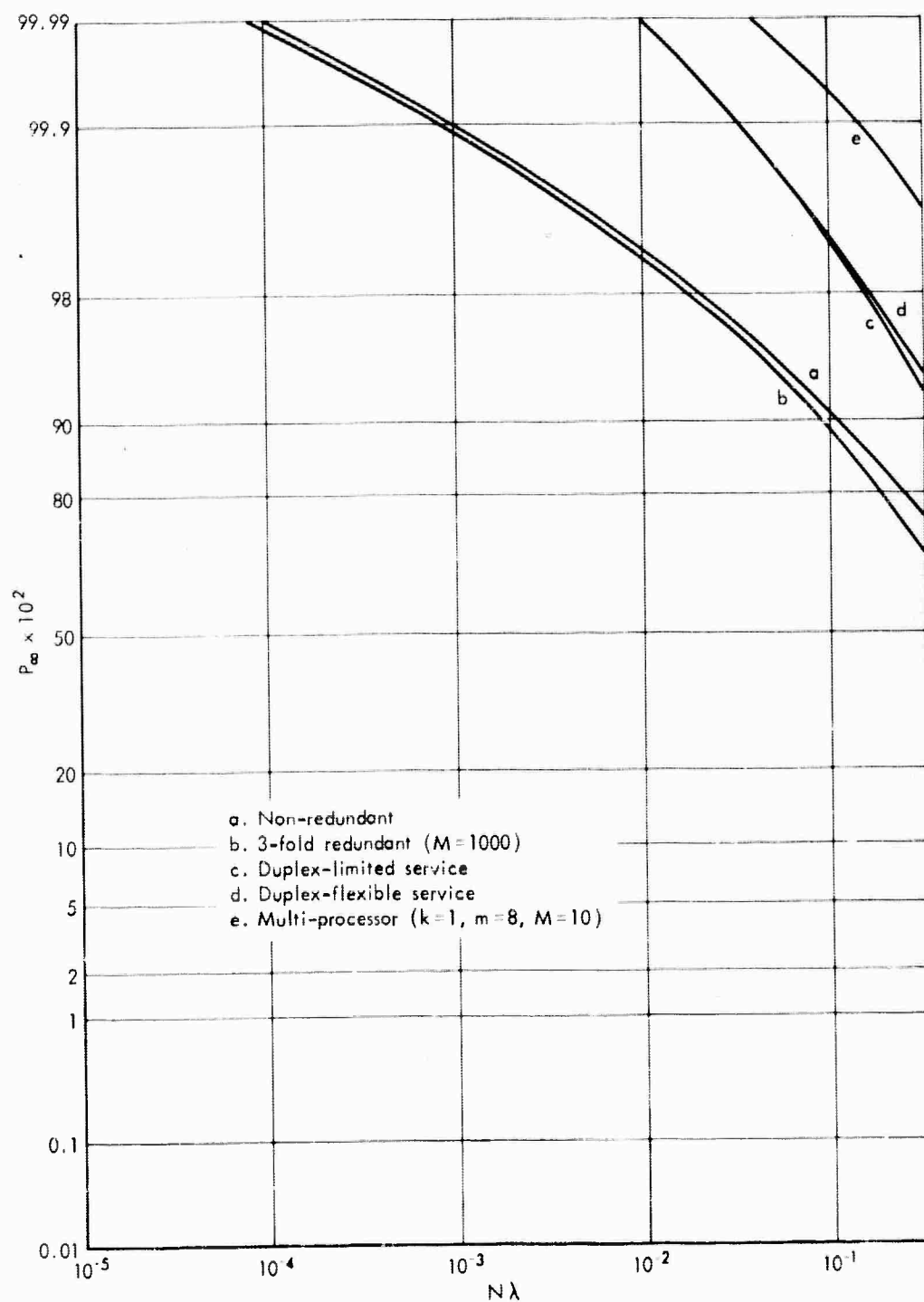


Fig. VII-3 — Comparison of systems ($\mu = 1.0$)

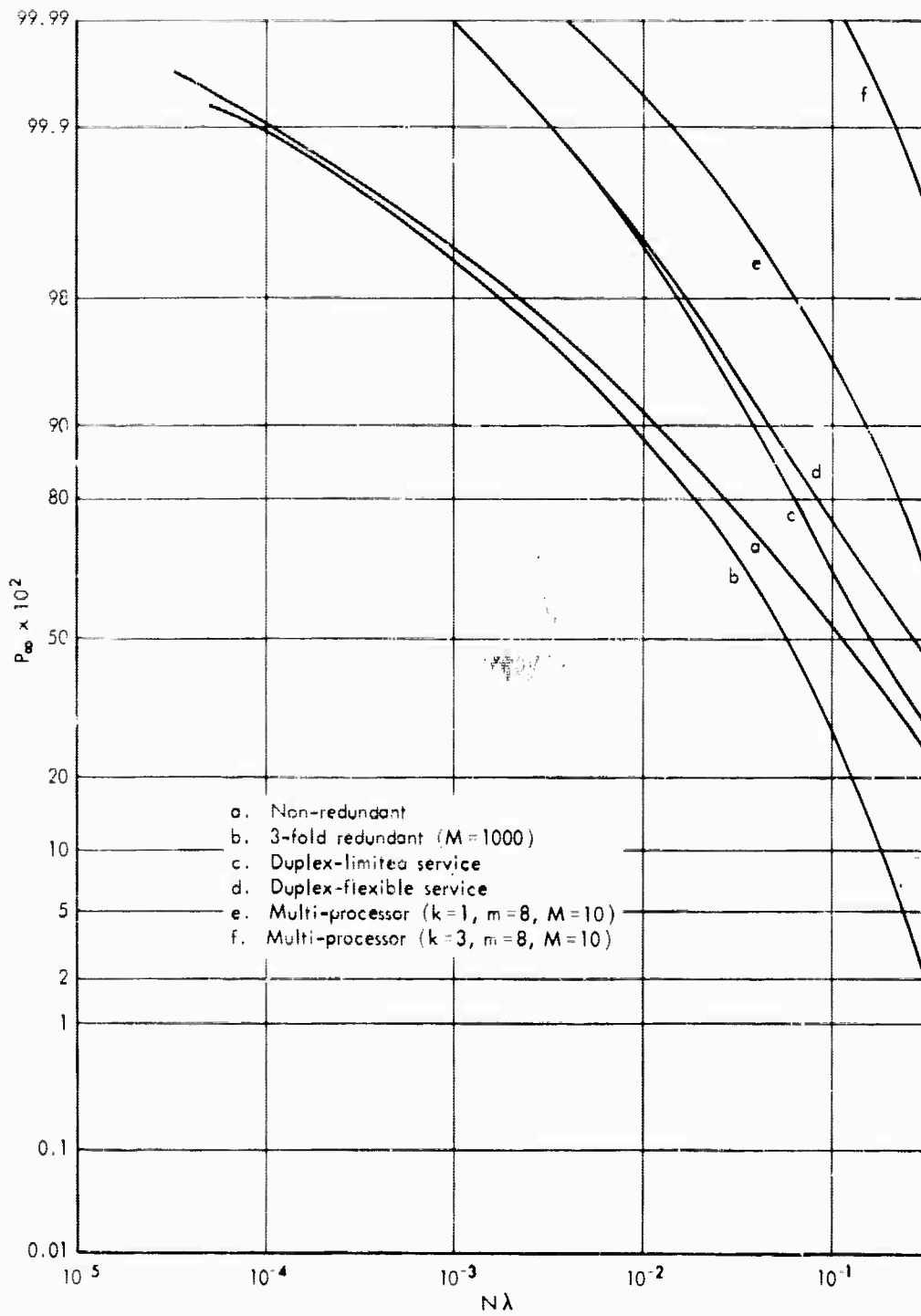


Fig. VII-4 — Comparison of systems ($\mu = 0.1$)

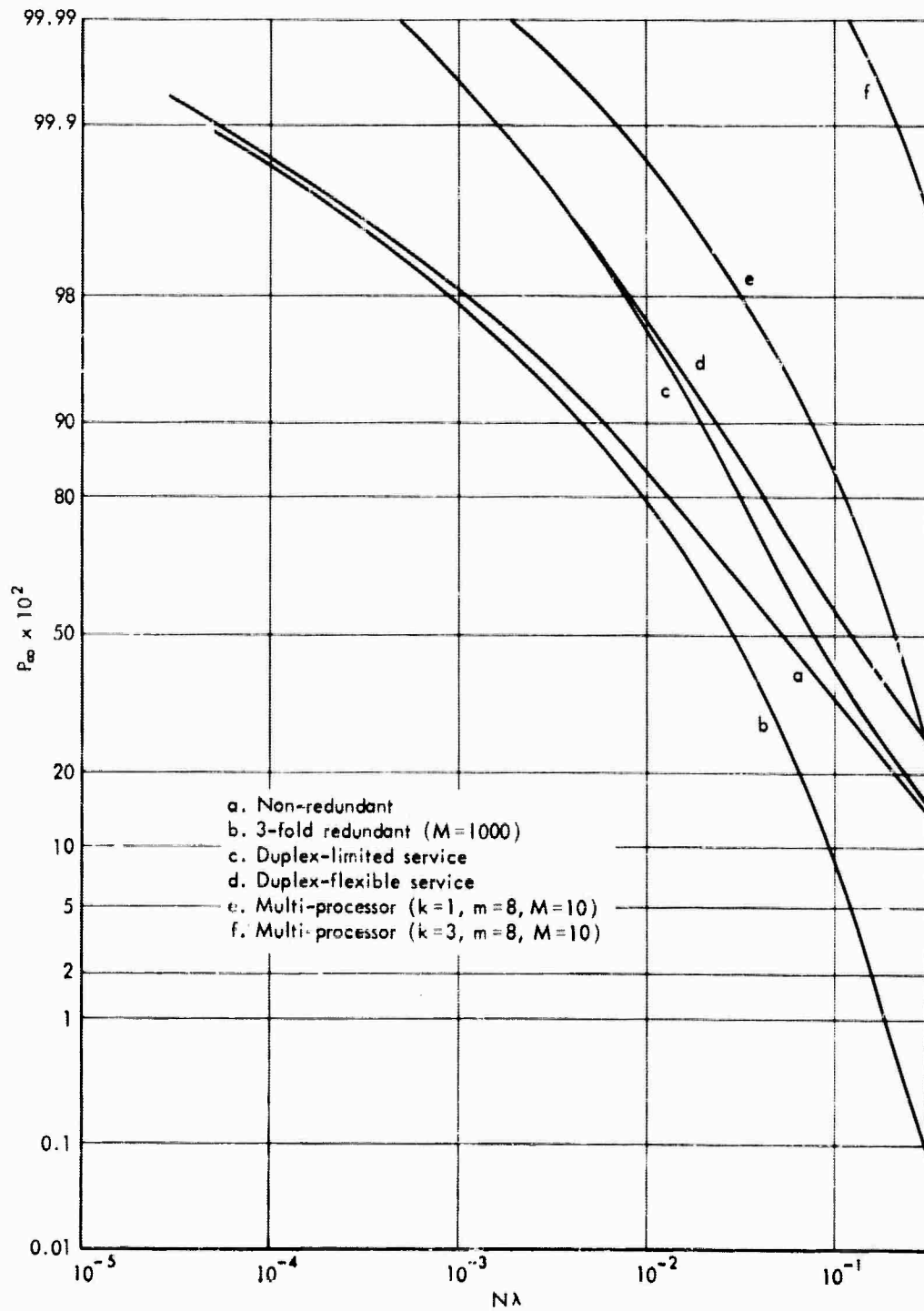


Fig.VII-5—Comparison of systems ($\mu=0.5$)

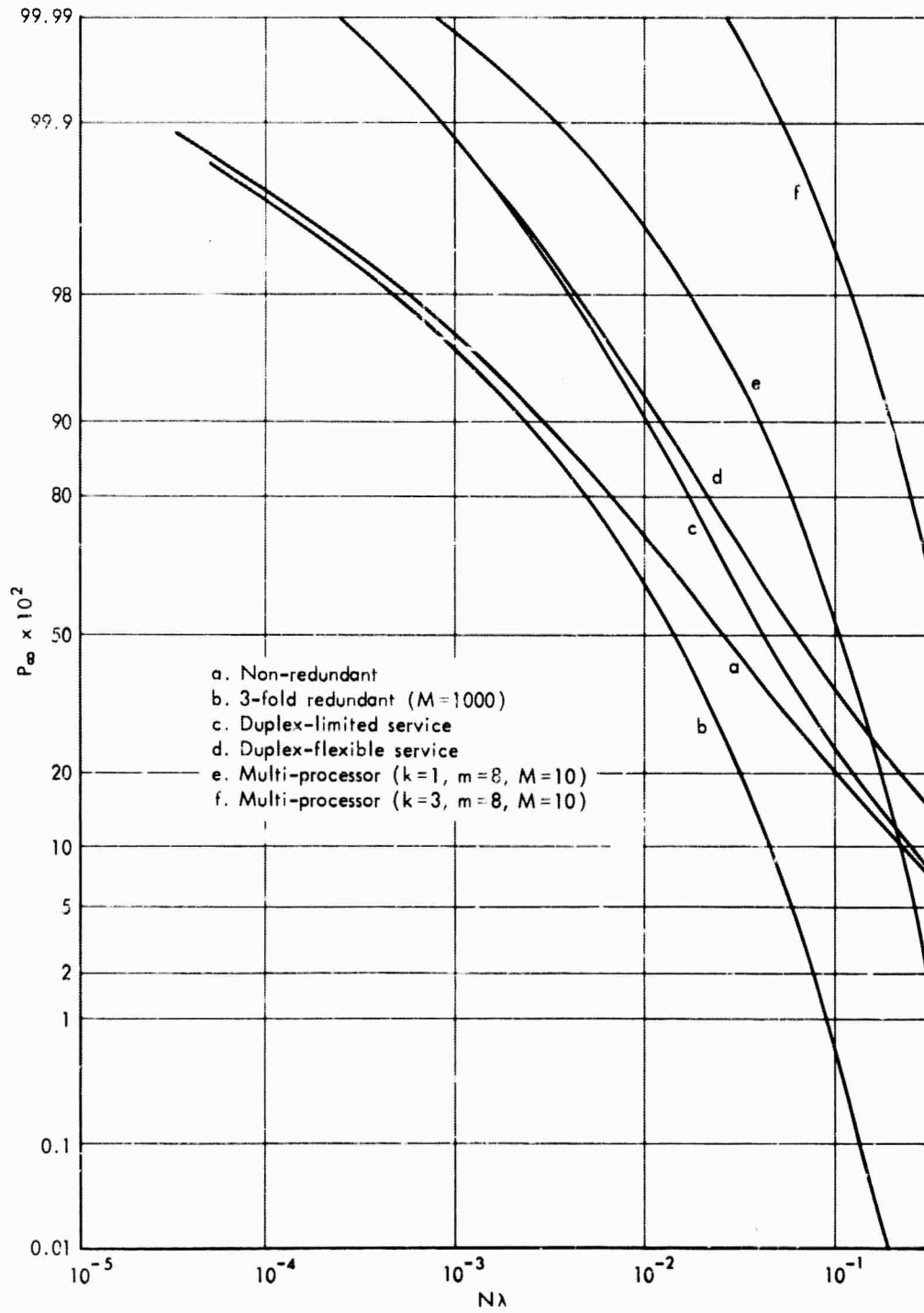


Fig.VII-6—Comparison of systems ($\mu = .025$)

not like "correct" circuit designs; they are either right or wrong and not, in theory, a matter of judgment. But the problem still remains as to how to tell if they are right. The nature of the problem is outlined and a modus operandi that should be helpful in obtaining error-free programs is presented in Chap. V. It must be read in its uncondensed form. Like so many other good and seemingly reasonable sets of rules, these do not become difficult until they are applied; whether a group of programmers (and their supervisors) could be made to follow the rules is speculative.

It is further concluded that very little can be done in terms of programmed error detection. "Programmed error detection" means the execution of a set of instructions whose purpose it is to verify the correctness of the computer logic or another set of instructions. Test problems, whose solutions are known, approximate an automatic check of the logic, but no program which checks another program yet exists. Simulation is a possibility and there is still room for further research, but all in all this is not a very optimistic area.

14. CONNECTORS AND PACKAGING

Connectors and packaging do not constitute a threat to availability if the best of today's techniques are employed. How to treat extreme environment (e.g., excessive radiation) is just barely touched upon near the end of Chap. V.

15. MANUFACTURER

Finally, Chap. V considers the role of the manufacturer. For a subject of such importance, and one which is so non-quantifiable, very little has been said here. The subject of the manufacturer's attitude and what helps or hinders the procurement of high-quality, reliable equipment needs much more attention.

16. AUTOMATIC FAULT DIAGNOSIS AND ISOLATION

Chapter VI outlines rather completely the current approaches to automatic fault diagnosis and isolation. It shows the success of any particular scheme depends heavily on pre-planned hardware additions, the quality of the service people, and the amount of money one wishes to spend.

Automatic fault diagnosis and isolation programs can be written (e.g., IBM System/360). They are effective (most single errors locatable to a single module), and reasonably complete (the single error must be excluded from only ten per cent of the machine).

17. OPTIMUM MODULE SIZE

Some analysis of the optimal size of a replaceable module concludes Chap. VI. This is considered by first estimating the relative time to diagnose the fault as a function of the module size, then estimating the relative cost of a set of spares as a function of module size. The two functions are summed with equal weights, and the resulting function is observed to have an absolute minimum--which implies the probability of an optimum size.

Appendix A

INTRODUCTION TO THE MATHEMATICS OF AVAILABILITY

1. INTRODUCTION

This appendix presents a brief derivation, from first principles, of the availability results cited in Chaps. III and IV. Nothing done here is either new or novel. Appendix A is included in this report not only for the reader's convenience (to aid in attacking problems not explicitly discussed in the main text), but also to support the conclusions presented above. Knowledge of elementary probability theory is assumed, and the presentation will be as brief as possible. For greater detail (and considerably more elegance) the reader is referred to Cox [1].

Before proceeding, some comments on transient and asymptotic (steady-state) solutions are in order. Most of the pertinent results concern a probability function, $P(t)$, usually called the "availability"--it is the probability that a system is available for use at time t . A "system" is anything from a single transistor to an entire computing complex, depending on the context. In some instances, this function is very difficult to evaluate (for either analytic or numerical reasons which will be

discussed presently), and the asymptotic value, $P_{\infty} = \lim_{t \rightarrow \infty} P(t)$, will be computed instead.

Care must be taken that P_{∞} in fact answers the proper questions about availability, because the system might have a transient period that is significant relative to the equipment lifetime; i.e., at the time we are concerned about the value of $P(t)$, it may not, in truth, be close enough to P_{∞} . It is argued here that for systems whose reliability must be very high and where service is possible, the use of P_{∞} is perfectly acceptable. The following reasons support this view:

- o For all cases of interest, $P_{\infty} < P(t)$, hence P_{∞} is at least a lower bound on the availability;
- o Most electronic design procedures such as "worst case" or "almost worst case" introduce more "over-design" into the system than would using P_{∞} as a measure of availability instead of the true $P(t)$.

In addition to these reasons which only say, in effect, that no harm is done by using P_{∞} , there are some more definite results.

- o For non-redundant systems, $P(t)$ is of such a form that very small transient periods imply large P_{∞} , and conversely. Therefore, if a highly reliable

system is desired, the system will have a very short transient phase, and the error in neglecting the transient contribution to $P(t)$ is negligible.

- o If service is possible, then the redundant computer (even with service) is not better (i.e., does not have larger $P(t)$) than the multiple machine configuration except for very small times. Redundant techniques are always accompanied by long transient periods. Redundant computers are considerably more complex, for fixed $P(t)$, than multiple computers, and their use in ground-based, serviceable systems is doubtful. Of course, without service, the redundant machine is the proper candidate, since the transient term accounts for most of the reliability--thus their favored use in space applications.

In the work that follows, the service time[†] will always be assumed to be exponentially distributed. When discussing the repair of a computer, this assumption seems entirely reasonable--the only other likely candidate being a service time of fixed duration. These two cases are the same for the asymptotic situation; i.e., if the mean of the exponential service distribution is equal to some assumed constant

[†]"Service time" is a random variable whose sample value is the length of time the computer is inoperative following a failure.

service time, then the probability of the machine being on, i.e., working properly, for large t is the same regardless of which distribution is assumed. This is not true for the transient phase, but the analysis is complicated and will not be carried out here. Readers desiring more information on the constant service time case are referred to Saaty [2].

2. DEFINITIONS AND THE POISSON PROCESS

First, some definitions and notation are necessary. Let T_f be a random variable denoting the time of failure (exactly what has failed, a single part or an entire system, will be clear from the context). Assume T_f has a distribution function, $F(t) = \Pr[T_f \leq t]$. The "failure rate," $r(t)$, is defined by the conditional probability[†]

$$r(t) = \lim_{\Delta t \rightarrow 0} \frac{\Pr[t \leq T_f \leq t + \Delta t | T_f \geq t]}{\Delta t} \quad (1)$$

If T_f also has a probability density function, $f(t)$, then

$$F(t) = \int_0^t f(x) dx \quad (2)$$

[†]Another definition of $r(t)$ is that it gives the probability of "immediate" failure of an item, given that the item has age t .

$$\bar{F}(t) = 1 - F(t) = \int_t^{\infty} f(x) dx \quad (3)$$

and

$$\begin{aligned} r(t) &= \lim_{\Delta t \rightarrow 0} \frac{\Pr[t \leq T_f \leq t + \Delta t, T_f \geq t]}{\Delta t \Pr[T_f \geq t]} = \lim_{\Delta t \rightarrow 0} \frac{\Pr[t \leq T_f \leq t + \Delta t]}{\Delta t \Pr[T_f \geq t]} \\ &= \frac{f(t)}{[1 - F(t)]} = \frac{f(t)}{\bar{F}(t)}. \end{aligned} \quad (4)$$

By integrating (4) we get

$$\int_0^t r(x) dx = \int_0^t \frac{dF(x)}{1-F(x)} = -\ln[1 - F(t)]. \quad (5)$$

Hence

$$F(t) = 1 - \exp \left[- \int_0^t r(x) dx \right] \quad (6)$$

and

$$f(t) = r(t) \exp \left[- \int_0^t r(x) dx \right]. \quad (7)$$

Failure processes whose failure rate, $r(t)$, is a constant are used frequently enough to make it worthwhile to derive these processes from first principles. We first require the notion of a Poisson process which has the very attractive property that the future behavior of the process is independent of the past. That is, the probability of an event occurring in an interval of length Δ depends only on the length of the interval and not on when Δ occurs in the time history of the process. A sufficient specification for a stochastic process, $f(t)$, to be a Poisson process is given by the following postulate [3,4]:

The stochastic process $f(t)$ is a Poisson process if, for Δ sufficiently small, there exists a positive constant λ such that the probability of one event occurring in $(t, t+\Delta)$ is $\lambda\Delta + o(t)$ and the probability of more than one event occurring is $o(t)$.[†]

The exponential failure distribution (constant failure rate) may be derived from this postulate. Consider a single part whose failure behavior obeys the postulate and is instantaneously replaced with an identical part if a failure does occur. Define a random variable N_t which

[†] $o(x)$ denotes a function that has the property $\lim_{x \rightarrow 0} o(x)/x = 0$.

denotes the number of failures in $(0, t)$ and let $P_m(t) = \Pr[N_t = m]$. Use of the postulate gives the following difference equation:

$$P_m(t+\Delta) = P_m(t)(1 - \lambda\Delta) + P_{m-1}(t)\lambda\Delta. \quad (8)$$

Rearranging and passing to the limit gives the differential equation

$$\lim_{\Delta \rightarrow 0} \frac{P_m(t+\Delta) - P_m(t)}{\Delta} = P'_m(t) = \lambda[P_{m-1}(t) - P_m(t)].^\dagger \quad (9)$$

For $m = 0$, Eq. (9) reduces to

$$P'_0(t) + \lambda P_0(t) = 0. \quad (10)$$

Solving (10) with the initial condition $P_0(0) = 1$, and then solving (9) repeatedly, we find that $P_m(t)$ is given by the Poisson distribution

[†]Use of the prime, e.g., $P'(t)$, will always denote differentiation with respect to time.

$$P_m(t) = \Pr[N_t = m] = \frac{(\lambda t)^m e^{-\lambda t}}{m!} . \quad (11)$$

Using the definition given in Eq. (3) and assuming a Poisson process therefore gives

$$\tilde{F}(t) = \Pr[T_f > t] = P_0(t) = e^{-\lambda t} . \quad (12)$$

Then $F(t)$ is the exponential distribution function

$$F(t) = \Pr[T_f \leq t] = 1 - e^{-\lambda t} , \quad (13)$$

Thus

$$f(t) = \lambda e^{-\lambda t} \quad (14)$$

and from Eq. (4)

$$r(t) = f(t)/\tilde{F}(t) = \lambda . \quad (15)$$

Also, the expectation of N_t is $1/\lambda$, which is commonly called the mean time to failure (MTTF).

3. NO REDUNDANCY--NO SERVICE AND EXPONENTIAL FAILURE DISTRIBUTION

As a rough approximation of reality, assume that a computer consists of N identical parts. Further, assume that the part failures are statistically independent, and subject to failures which are exponentially distributed with failure rate λ .

Then the probability that a part survives past time t is given by

$$\bar{F}_p(t) = e^{-\lambda t} \quad (16)$$

and the probability that the entire computer survives past t , $\bar{F}_c(t)$, (given that it was on at $t=0$), is equal to the probability that all N parts have survived past t ,[†]

$$\bar{F}_c(t) = e^{-N\lambda t} . \quad (17)$$

For completeness more than any other reason, $\bar{F}_c(t)$ is shown in Fig. A-1 for some interesting values of $N\lambda$.

[†]This assumption, that the failure of any part causes the failure of the machine, is the standard approach. However, in practice it might not be the case, since it is dependent on the nature of the program.

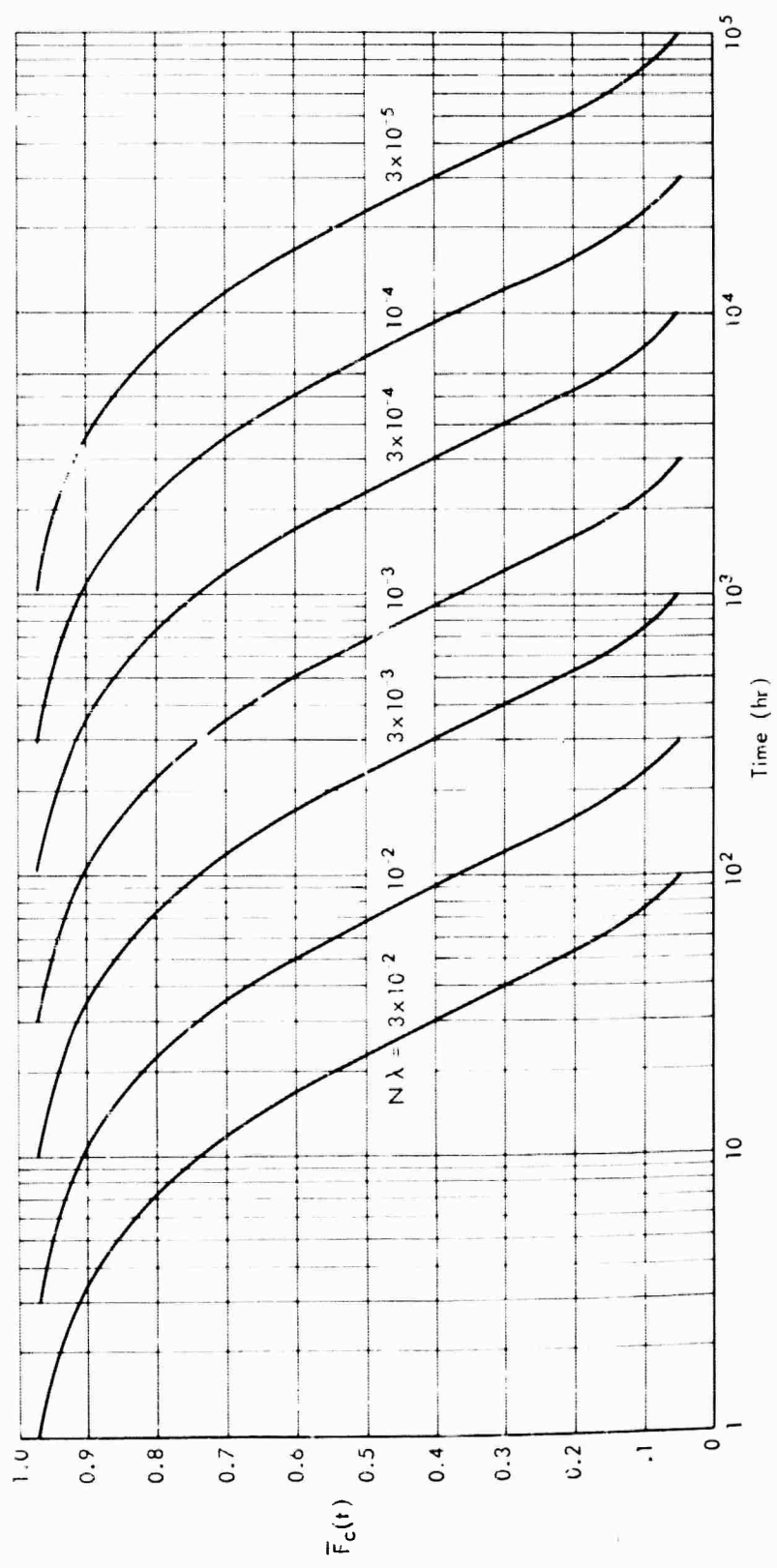


Fig.A-1 —The exponential distribution function, $\bar{F}_c(t) = e^{-N\lambda t}$

For a slightly more accurate model, assume that a computer consists of N_1 transistors, N_2 diodes, and N_3 resistors and capacitors, where $N_1 + N_2 + N_3 = N$. If a transistor has failure rate λ_1 , diodes λ_2 , and resistors and capacitors λ_3 , then the probability that the computer is operating after time t is $F_c(t) = \exp[-(N_1\lambda_1 + N_2\lambda_2 + N_3\lambda_3)t]$.

4. 2m+1-FOLD REDUNDANCY WITH PERFECT VOTING--NO SERVICE AND EXPONENTIAL FAILURE DISTRIBUTION

The purpose of this and the next few sections is to give a general probabilistic description of machines which employ redundancy to achieve reliability. Knox-Seith [5] and Wilcox and Mann [6] give a much more detailed analysis of the problem.

First consider a collection of M subsystems, $1 \leq M \leq N$, which might constitute an N -part computer (N/M is an integer). The size of M defines the level at which the redundancy will be employed. For statistical purposes, all M subsystems will be assumed identical. For instance, if $M = N$, then each individual part is a "subsystem"; but if $M = 1$, the entire computer is the subsystem. The intent is to reproduce each of the M subsystems $2m+1$ times in order to gain the benefits of a redundant system. This is illustrated in Fig. A-2.

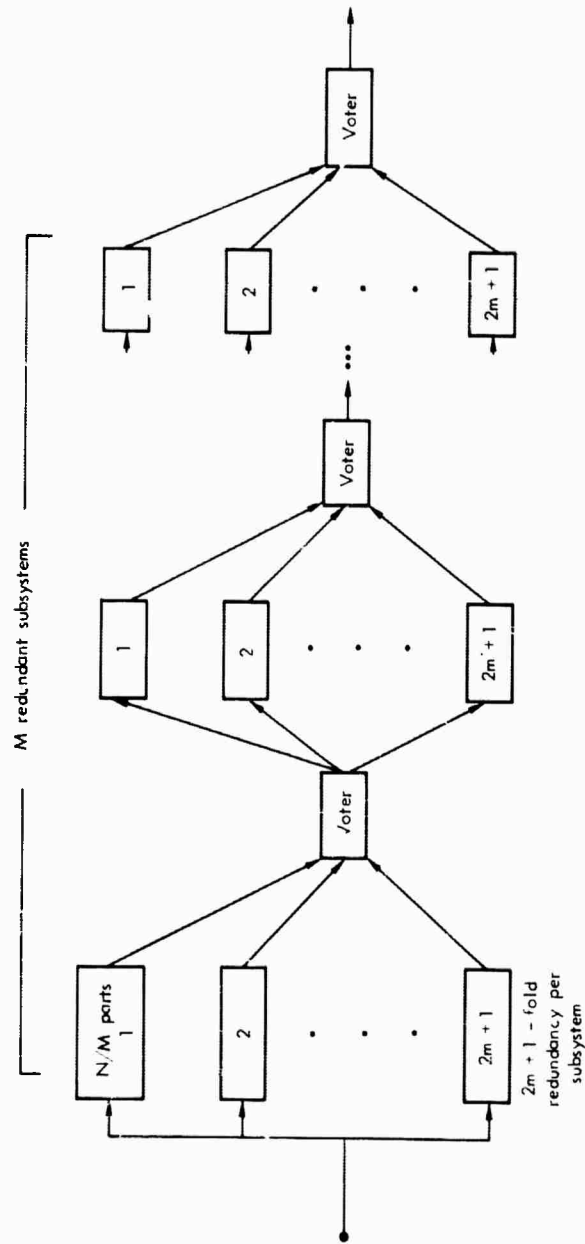


Fig.A-2—Redundant computer with majority voting

Each subsystem contains N/M parts and all parts are independent and have identical exponential failure distributions. The probability of survival of any non-redundant subsystem is

$$\bar{F}_s(t) = \exp(-N\lambda t/M) . \quad (18)$$

Now with $2m+1$ -fold redundancy and failure-free majority voting, a subsystem survives if at least $m+1$ of the circuits making up the redundant group survive. The probability of this event is given by

$$\begin{aligned} P_s(t) &= \Pr[k \geq m+1] = \sum_{k=m+1}^{2m+1} b[k; 2m+1, \bar{F}_s(t)] \\ &= 1 - \sum_{k=0}^m b[k; 2m+1, \bar{F}_s(t)] \end{aligned} \quad (19)$$

where k is the number of circuits which survive past t and $b(k; n, p) = n! p^k (1-p)^{n-k} / k! (n-k)!$. The failure of any redundant group can cause the failure of the computer. Hence the probability that the computer survives until t , assuming independent subsystems, is

$$P_c(t) = [P_s(t)]^M$$

$$= \left[1 - \sum_{k=0}^m \frac{(2m+1)!}{k!(2m+1-k)!} \left(e^{-\frac{N\lambda t}{M}} \right)^k \left(1 - e^{-\frac{N\lambda t}{M}} \right)^{2m+1-k} \right]^M$$

(20)

$P_c(t)$ is shown graphically in Figs. A-3 to A-10 for $m=1,2$ and some typical values of $N\lambda$ and M . If the assumption of perfect voting is used, inspection of Eq. (20) shows that $P_c(t)$ is maximized (for fixed N and λ) by choosing M as small as possible; i.e., the voting should be done at the lowest possible level.

5. N-FOLD REDUNDANCY WITH IMPERFECT VOTING--NO SERVICE AND EXPONENTIAL FAILURE DISTRIBUTION

If the voters themselves are also susceptible to failure, then redundant voters should be used to increase the voter reliability. This results in the system shown in Fig. A-11. Again, each non-redundant subsystem has N/M parts and each subsystem has independent and identical survival probability, $P_s(t) = \exp(-N\lambda t/M)$. When using redundant voters, it is not necessarily optimum to vote each redundant subgroup. Instead, L groups are connected in series prior to each voting operation. The probability

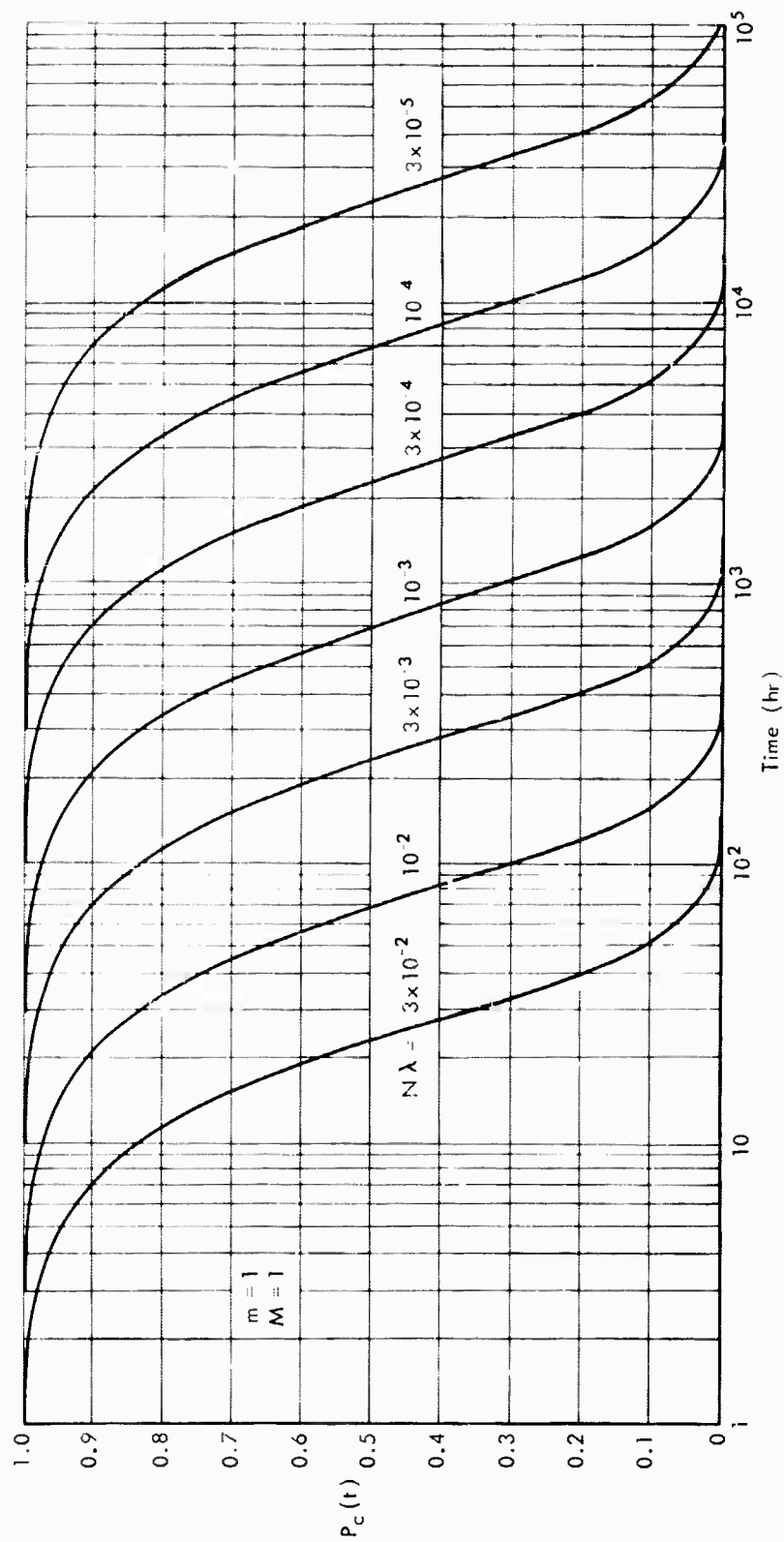


Fig. A-3 ---Availability of redundant computer---no service and exponential failure

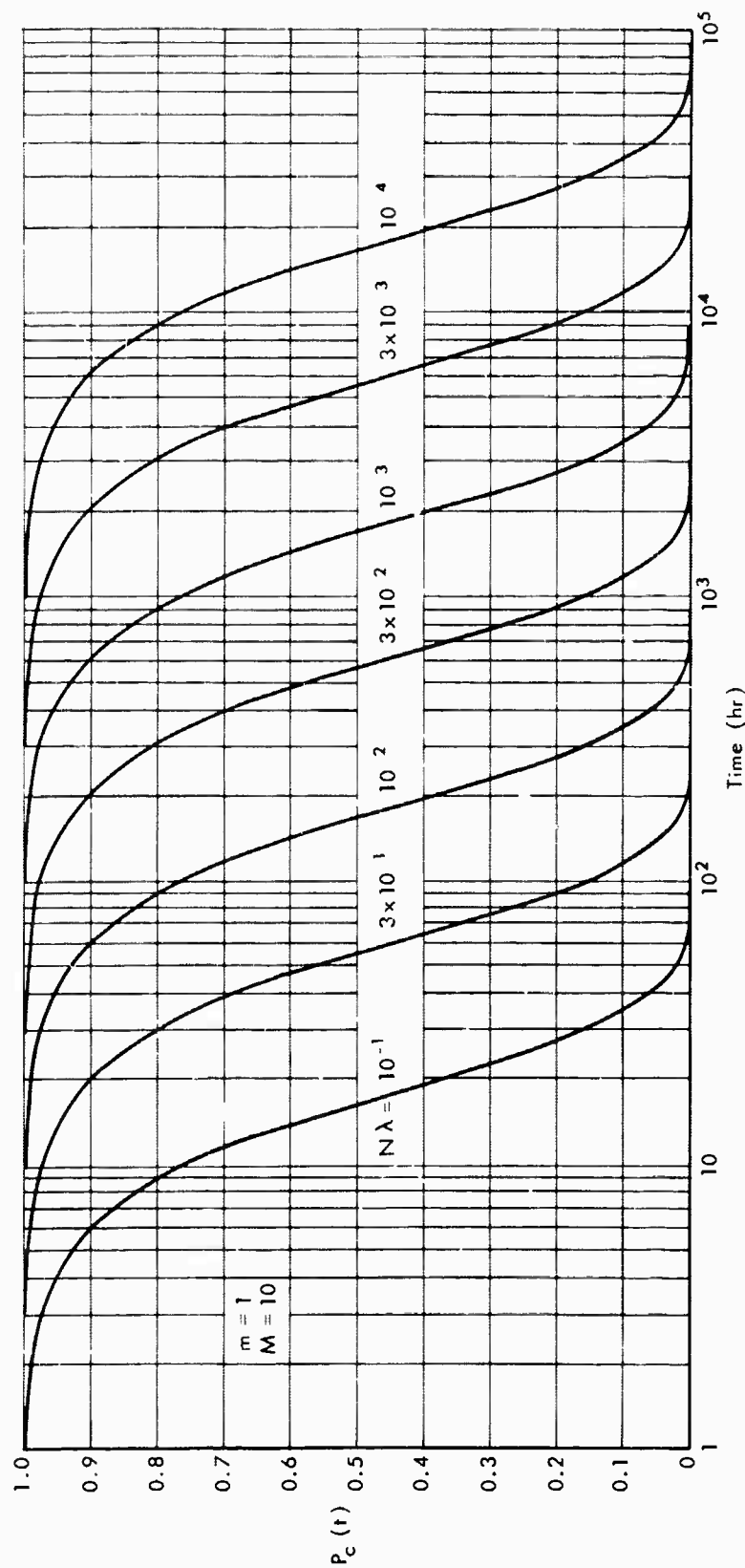


Fig. A-4 —Availability of redundant computer—no service and exponential failure

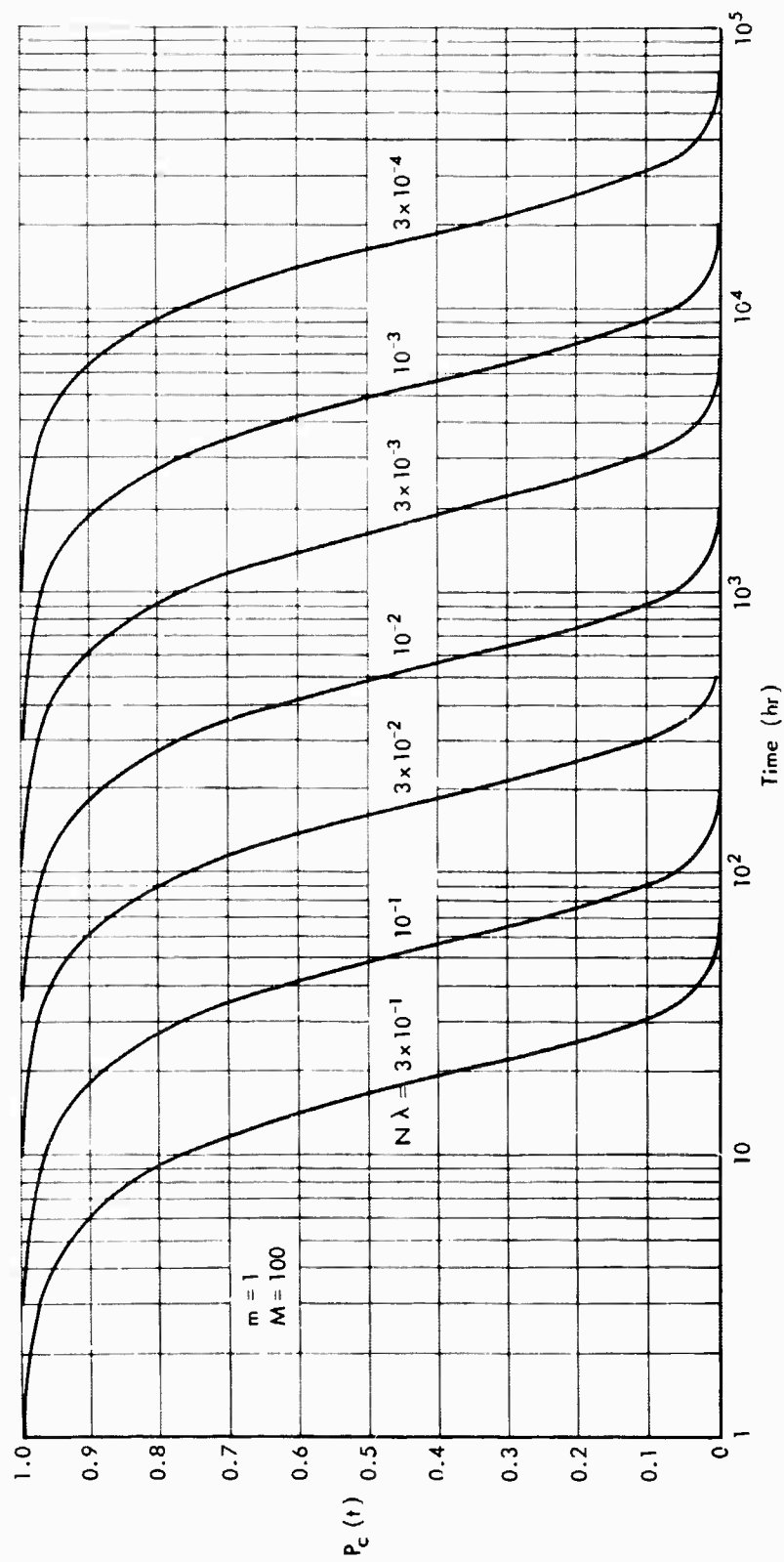


Fig.A-5—Availability of redundant computer—no service and exponential failure

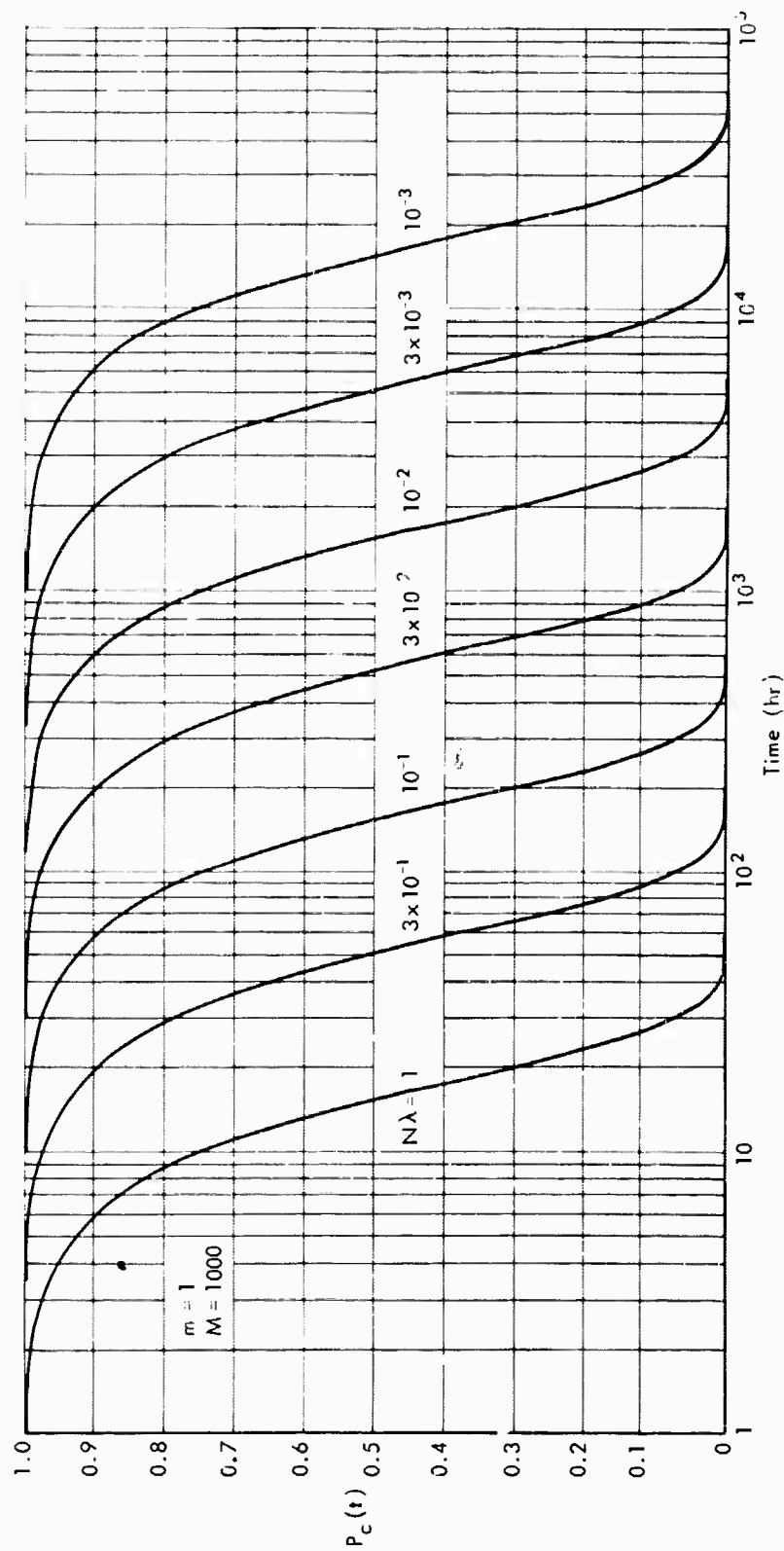


Fig. A-6 --- Availability of redundant computer---no service and exponential failure

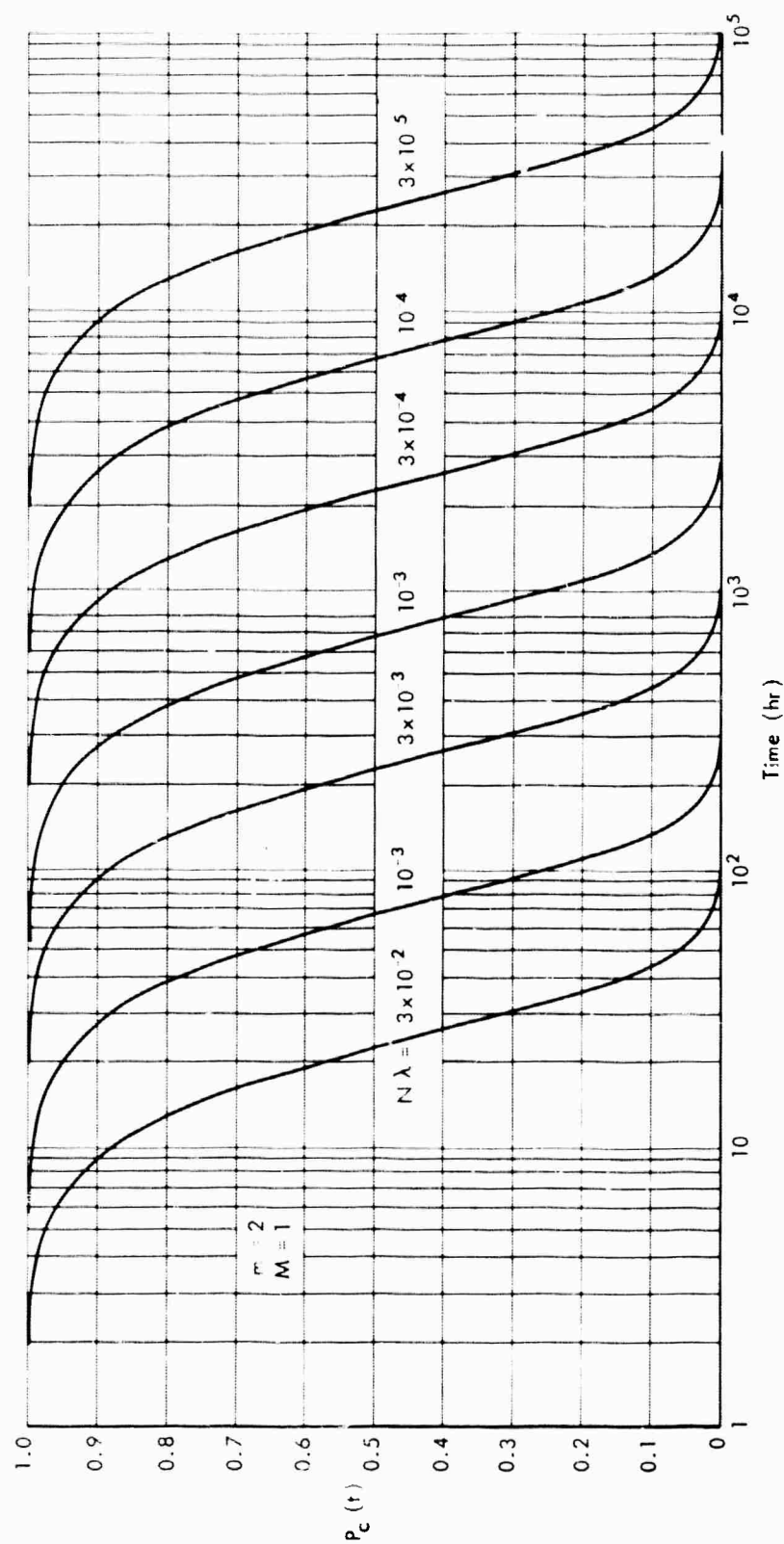


Fig.A-7 — Availability of redundant computer—no service and exponential failure

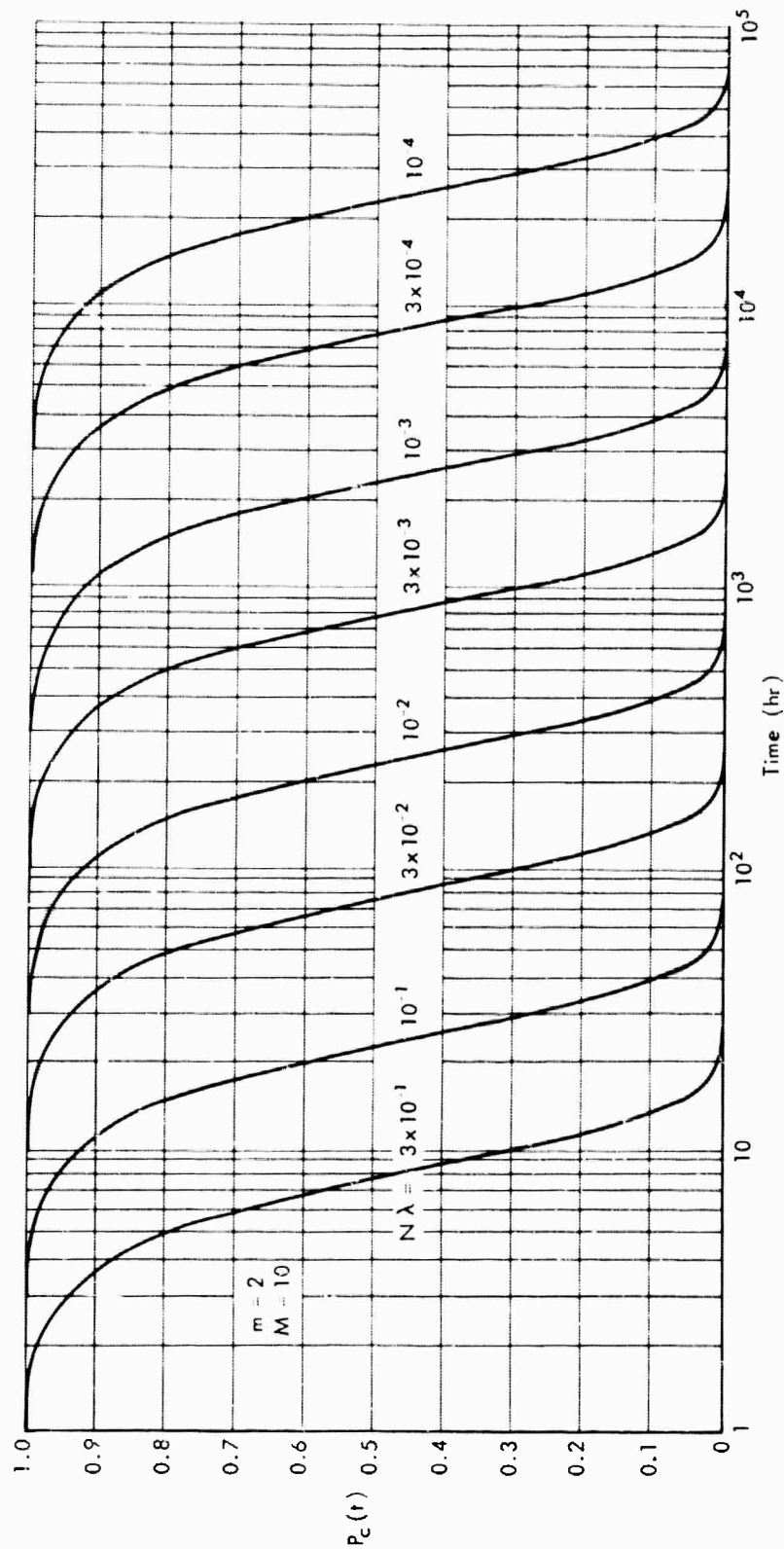


Fig. A-8—Availability of redundant computer—no service and exponential failure

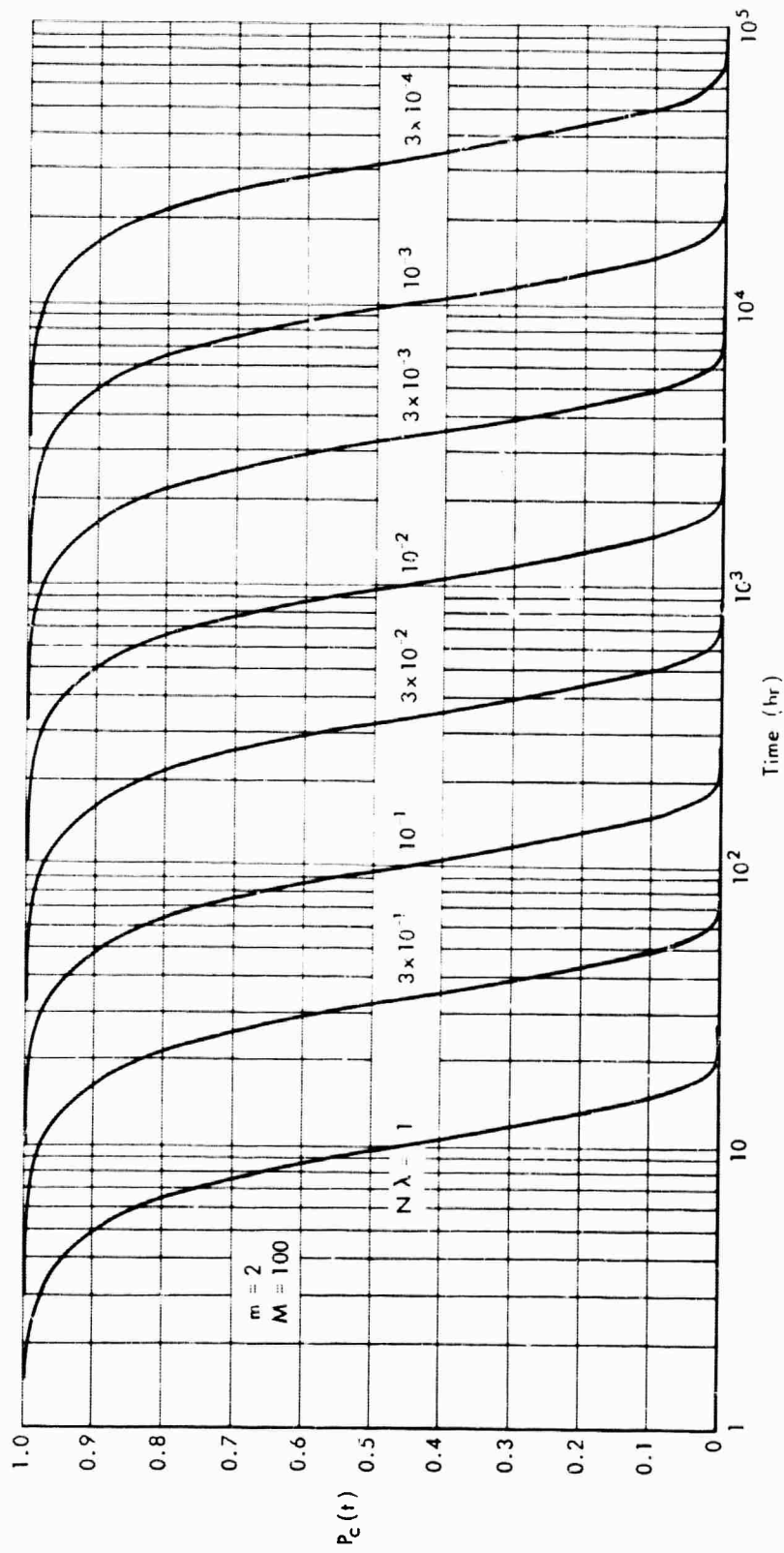


Fig. A-9—Availability of redundant computer—no service and exponential failure

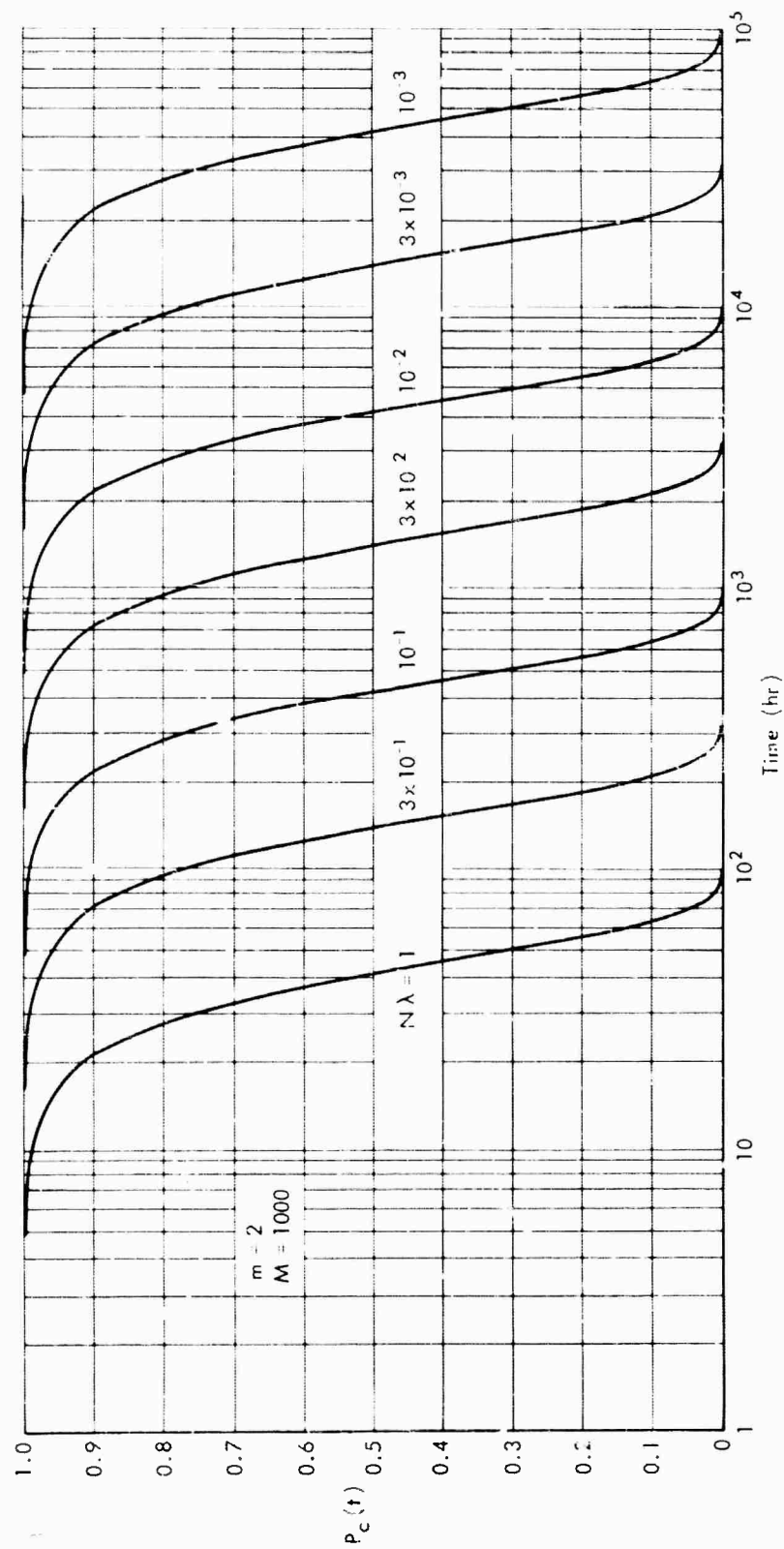


Fig.A-10—Availability of redundant computer—no service and exponential failure

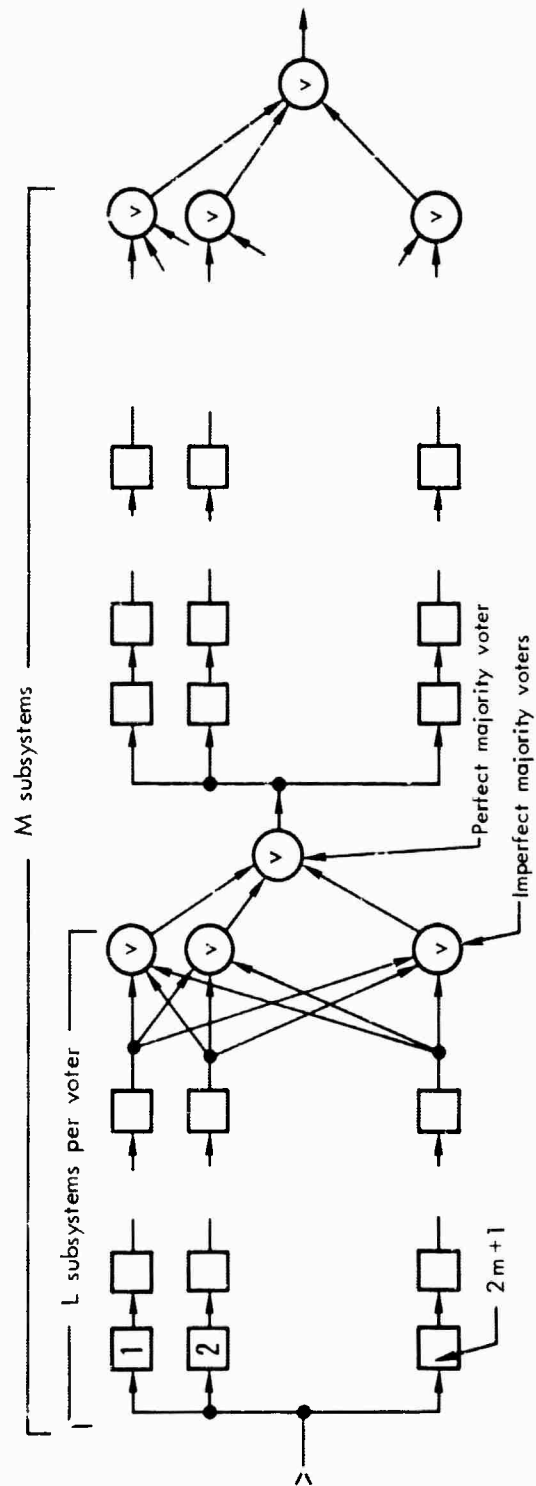


Fig.A-11 — Redundant computer with imperfect voters

that a series chain of L non-redundant subsystems survives until t is $P_L(t) = [P_s(t)]^L = \exp(-LN\lambda t/M)$. Now with $2m+1$ redundancy, the output of which will be majority-voted, the probability, $P_r(t)$, that there are at least $m+1$ L -chains is

$$P_r(t) = \sum_{k=m+1}^{2m+1} b[k; 2m+1, P_L(t)] . \quad (21)$$

Next, assume that every voter is independent and has survival probability $\bar{F}_v(t) = \exp(-\xi t)$. At least $m+1$ must be operating for every $2m+1$ redundant chain. The survival probability, $P_{rv}(t)$, for the redundant L -chain plus voter is

$$P_{rv}(t) = \sum_{k=m+1}^{2m+1} \sum_{j=m+1}^{2m+1} b[k; 2m+1, P_L(t)] b[j; 2m+1, \bar{F}_v(t)] . \quad (22)$$

Finally, the probability that the entire computer survives is

$$P_c(t) = [P_{rv}(t)]^{M/L} . \quad (23)$$

No numerical examples of Eq. (23) are presented here, since only perfect voting is considered in the text. It is interesting to note before leaving this topic that with

the added parameter, ξ , contributed by the voter survival probability, there is a truly optimum value of $L > 1$. It is no longer true that voting should be done at the lowest level, but L is a function of λ and ξ .

6. NO REDUNDANCY--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTION

By arguments given previously, we have shown that T_f , the time to first failure of an entire non-redundant N -component computer, is exponentially distributed with failure rate, $r(t) = N$ $f(t) = N\lambda \exp(-N\lambda t)$, and $\bar{F}(t) = \exp(-N\lambda t)$ (if, of course, part failures are exponentially distributed). Now define a new random variable, T_s , called the service time (the time to repair the computer), and take this also to be exponentially distributed with service rate μ . The pertinent distributions for the service time are

$$G(t) = \Pr[T_s \leq t] = 1 - e^{-\mu t}, \quad (24)$$

$$\bar{G}(t) = \Pr[T_s > t] = e^{-\mu t}, \quad (25)$$

and

$$g(t) = G'(t) = \mu e^{-\mu t}. \quad (26)$$

The question now is what is the probability, $P(t)$, that the computer is on at time t ?

Remembering that both the failure and repair distributions are derived from Poisson processes, and hence have no memory, we may write a different equation for $P(t)$. Namely, for Δ sufficiently small

$$P(t+\Delta) = P(t)\bar{F}(\Delta) + [1 - P(t)]G(\Delta) \quad (27)$$

where $\bar{F}(\Delta) = \exp(-N\lambda\Delta)$ and $G(\Delta) = 1 - \exp(-\mu\Delta)$. In other words, the probability that the computer is on at $t+\Delta$ is equal to the probability that the computer is on at t and remains on for a small additional time Δ , or that the computer was off at time t and repair was completed in an additional time Δ . The value of Δ is chosen to be so small that the probability of more than one repair or failure is negligible. Expanding $\bar{F}(\Delta)$ and $G(\Delta)$ in a Taylor's series and retaining only the first-order terms transforms Eq. (27) to

$$P(t+\Delta) = P(t)(1 - N\lambda\Delta) + [1 - P(t)]\mu\Delta. \quad (28)$$

Rearranging and taking the limit gives

$$\lim_{\Delta \rightarrow 0} \frac{P(t+\Delta) - P(t)}{\Delta} = P'(t) = - (N\lambda + \mu)P(t) + \mu. \quad (29)$$

With the choice of initial condition, $P(0)$, Eq. (29) is easily solved for $P(t)$. The most interesting initial condition is $P(0) = 1$. We find that the probability that the computer is on at any time t is

$$P(t) = \frac{\mu}{N\lambda + \mu} + \frac{N\lambda}{N\lambda + \mu} e^{-(N\lambda + \mu)t} . \quad (30)$$

$P(t)$ is usually called the "availability" (or sometimes, the "readiness"). $P_{\infty} = \lim_{t \rightarrow \infty} P(t) = \mu / (N\lambda + \mu)$. Figure A-12 shows P_{∞} for selected values of $N\lambda$ and μ .

7. REDUNDANT COMPUTER--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

If redundancy is now introduced into the computer at some subsystem level, and a constant service rate is maintained, the problem becomes considerably more complex. In the following analysis, the case of perfect voting will always be assumed.

Let $P_s(t)$ be the probability that a single redundant subsystem is on at time t . Next define two random variables: T_f , the length of time the subsystem is on before its first failure starting at $t=0$; and T_s , the length of time the subsystem is off immediately following T_f . As usual, $\Pr[T_f \leq t] = F(t)$ and $\Pr[T_s \leq t] = G(t)$. Now define

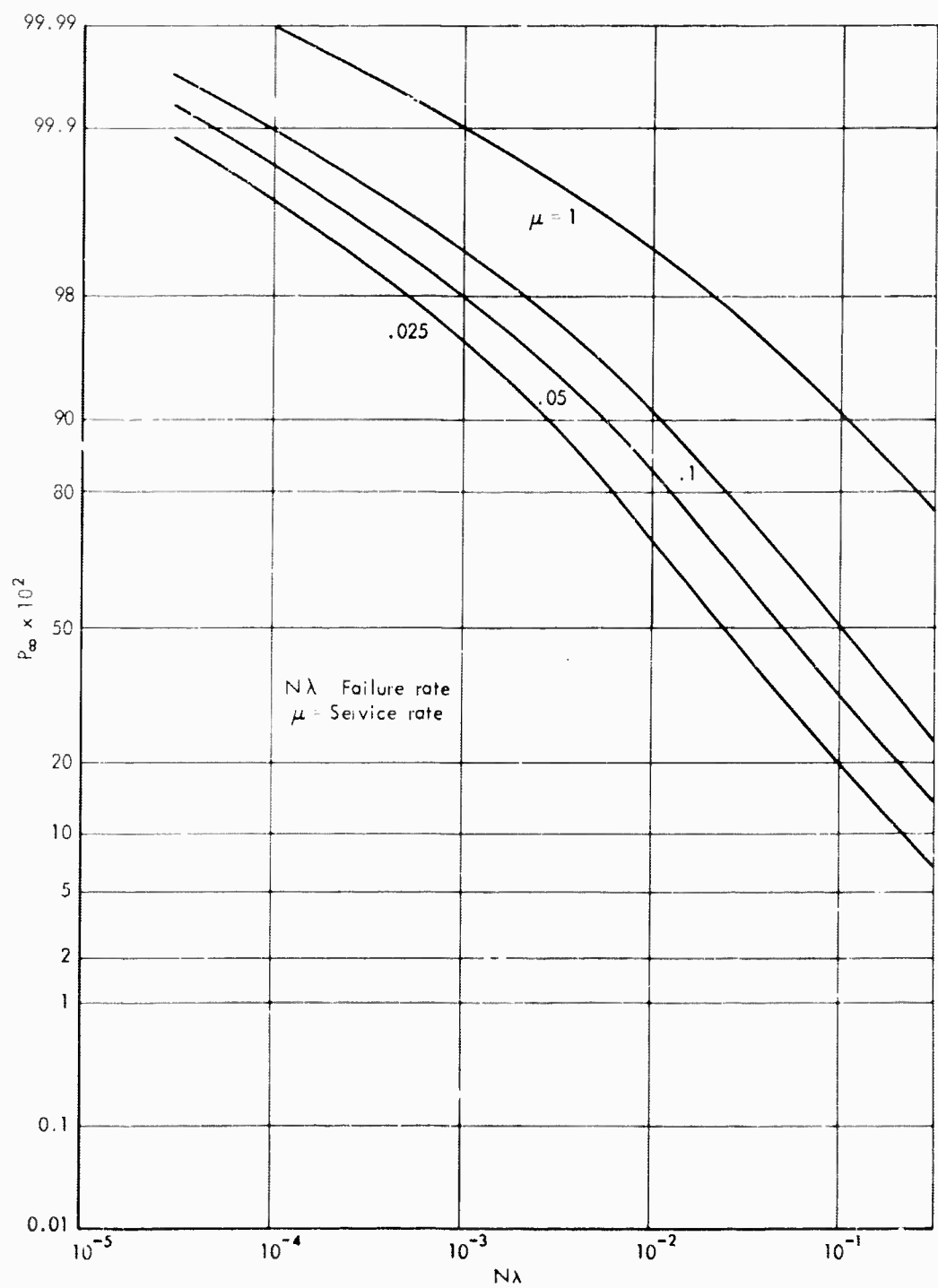


Fig.A-12—Asymptotic availability of non-redundant computer
(exponential failure and service)

a new random variable $T = T_f + T_s$, where T is the time at which the subsystem is on again after its first failure, and let $\Pr[T \leq t] = H(t)$. We know that

$$P_s(t) = \int_0^{\infty} P_s(t|T=u) dH(u) . \quad (31)$$

If $u \leq t$,

$$P_s(t|T=u) = P_s(t-u) . \quad (32)$$

Equation (32) is valid only if the entire subsystem is replaced at $t = u$, thus insuring a renewal process, or if the failure process is a Markov process (a class of stochastic processes of which the Poisson processes are a subset), i.e., the future state of each subsystem is dependent only on the present state and not on the past.

If $u > t$,

$$P_s(t|T=u) = P_s(t < T_f|T=u) \quad (33)$$

and

$$\begin{aligned} \int_0^{\infty} P_s(t|T=u) dH(u) &= \int_t^{\infty} P_s(t < T_f|T=u) dH(u) \\ &= P_s(t < T_f) = 1 - F(t) . \end{aligned} \quad (34)$$

Substituting Eqs. (32) and (34) into Eq. (31) gives the following integral equation:

$$P_s(t) = 1 - F(t) + \int_0^t P_s(t-u) dH(u) . \quad (35)$$

Assuming that T_f , T_s , and T have density functions f , g , and h respectively, and if T_f and T_s are independent, then

$$h(T) = \int_0^T f(T-v)g(v)dv \quad (36)$$

and

$$P_s(t) = 1 - F(t) + \int_0^t P_s(t-u)h(u)du . \quad (37)$$

Equation (37) is a Volterra integral equation of the second kind and its solution will provide the desired $P_s(t)$.

One method is to take the Laplace transform of Eq. (37):[†]

[†]The one-sided Laplace transform, $f^*(s)$, of $f(t)$, is defined as $f^*(s) = \int_0^\infty e^{-st}f(t)dt$. Operations such as taking

the transform of an integral or of a convolution are performed without explanation and the reader should consult any text on the subject (such as Gardner and Barnes [7]).

$$P_s^*(s) = \frac{1}{s} - \frac{f^*(s)}{s} + P_s^*(s)f^*(s)g^*(s) . \quad (38)$$

Rearranging gives the desired transform of $P_s(t)$:

$$P_s^*(s) = \frac{1 - f^*(s)}{s[1 - f^*(s)g^*(s)]} . \quad (39)$$

So far, our attention has been directed toward obtaining the entire transient solution to the availability problem, namely $P_s(t)$.[†] Success in this endeavor depends on a) the ability to obtain the transform $f^*(s)$ and $g^*(s)$, and b) the inversion of $P_s^*(s)$. Only in a few relatively simple examples does it appear possible to get the solution by the method of Laplace transforms. Such an example will be given next. Theorems do exist which relate P_∞ directly to $f(T_f)$ and $g(T_s)$ and will be used when required.

To continue with the case at hand, from Eq. (19) (with a slight shift in notation), we have

$$F(t) = \sum_{k=0}^m \frac{(2m+1)!}{k!(2m+1-k)!} \left(e^{\frac{-N\lambda t}{M}} \right)^k \left(1 - e^{\frac{-N\lambda t}{M}} \right)^{2m+1-k} \quad (40)$$

[†]Actually, we want $[P_s(t)]^M$, assuming M independent and identical subsystems.

and

$$G(t) = 1 - e^{-\mu t} . \quad (41)$$

Specializing to the case $m=1$ (three-fold redundancy) and evaluating Eq. (40) gives

$$F(t) = 1 - 3e^{-2\eta t} + 2e^{-3\eta t} , \quad \eta = \frac{N}{M}\lambda . \quad (42)$$

Differentiating Eqs. (42) and (41) results in

$$f(t) = 6\eta \left(e^{-2\eta t} - e^{-3\eta t} \right) \quad (43)$$

and

$$g(t) = \mu e^{-\mu t} , \quad (44)$$

respectively. These densities have the following transforms:

$$f^*(s) = \frac{6\eta^2}{(s+2\eta)(s+3\eta)} \quad (45)$$

and

$$g^*(s) = \frac{\mu}{s+\mu} . \quad (46)$$

Substituting Eqs. (45) and (46) into (39), we have

$$P_s^*(s) = \frac{s^2 + (5\eta + \mu)s + 5\mu\eta}{s[s^2 + (5\eta + \mu)s + 6\eta^2 + 5\mu\eta]} . \quad (47)$$

The inverse, $P_s(t)$, may be found by first evaluating the poles of $P_s^*(s)$, then using transform 1.111 of Ref. 7 if the poles are real or transform 5.2-(a) of Ref. 8 if the poles are complex. This work will not be shown here, but the results are shown for the selected cases $M=1$, $N\lambda=10^{-4}$, 10^{-3} , 10^{-1} in Figs. A-13 to A-15.

Proceeding next to the asymptotic situation, we wish to show that if T_f and T_s are independent then

$$\lim_{t \rightarrow \infty} P_s(t) = \frac{E[T_f]}{E[T_f] + E[T_s]} . \quad (48)$$

The final value theorem states that if $sP_s^*(s)$ is analytic on the axis of imaginaries and in the right half-plane [7], then

$$P_\infty = \lim_{t \rightarrow \infty} P_s(t) = \lim_{s \rightarrow 0} sP_s^*(s) . \quad (49)$$

Therefore, from Eq. (39),

$$P_\infty = \lim_{s \rightarrow 0} \frac{1 - f^*(s)}{1 - f^*(s)g^*(s)} \quad (50)$$

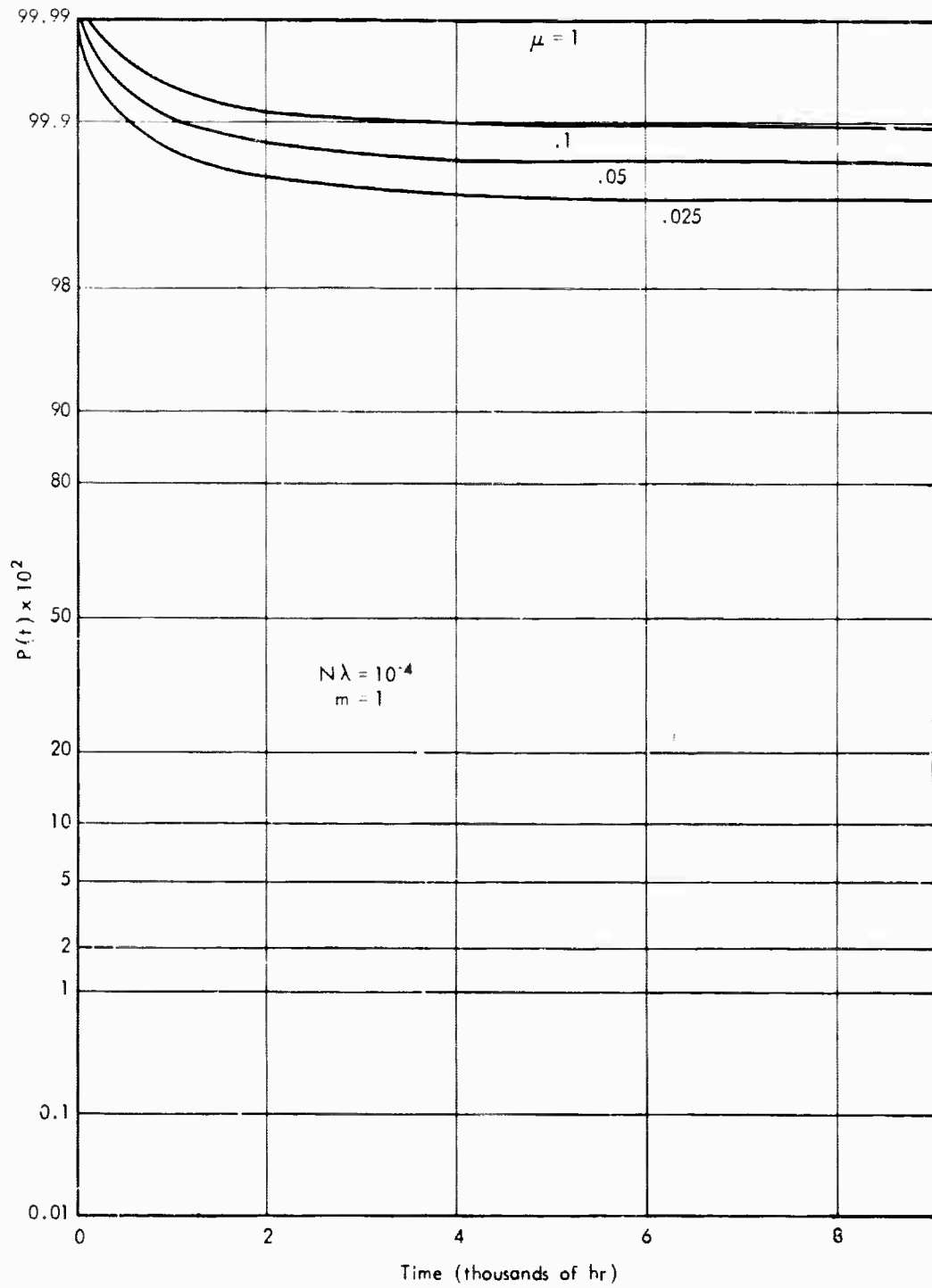


Fig.A-13—Availability of redundant computer
(transient phase, exponential service)

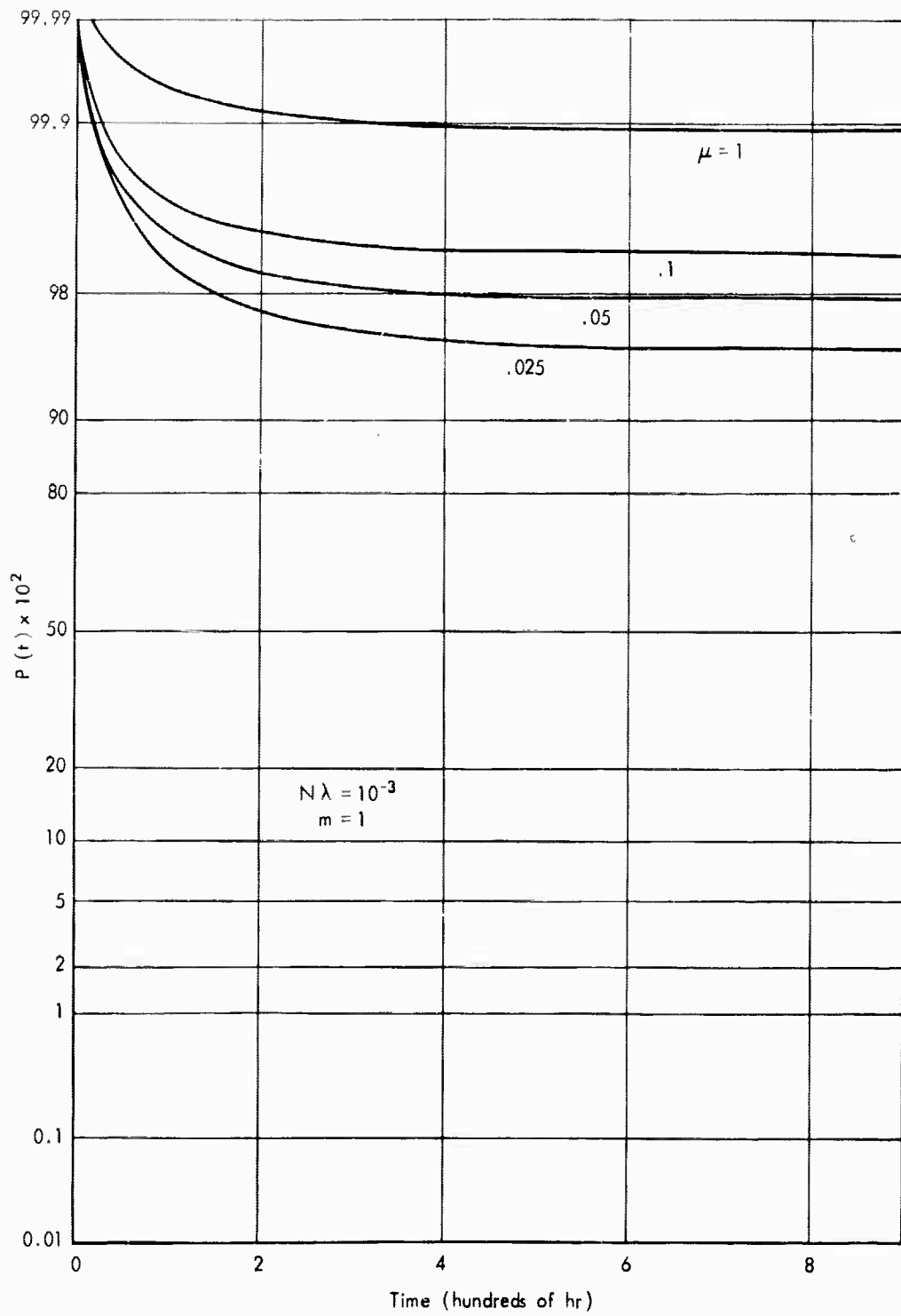


Fig.A-14—Availability of redundant computer
(transient phase, exponential service)

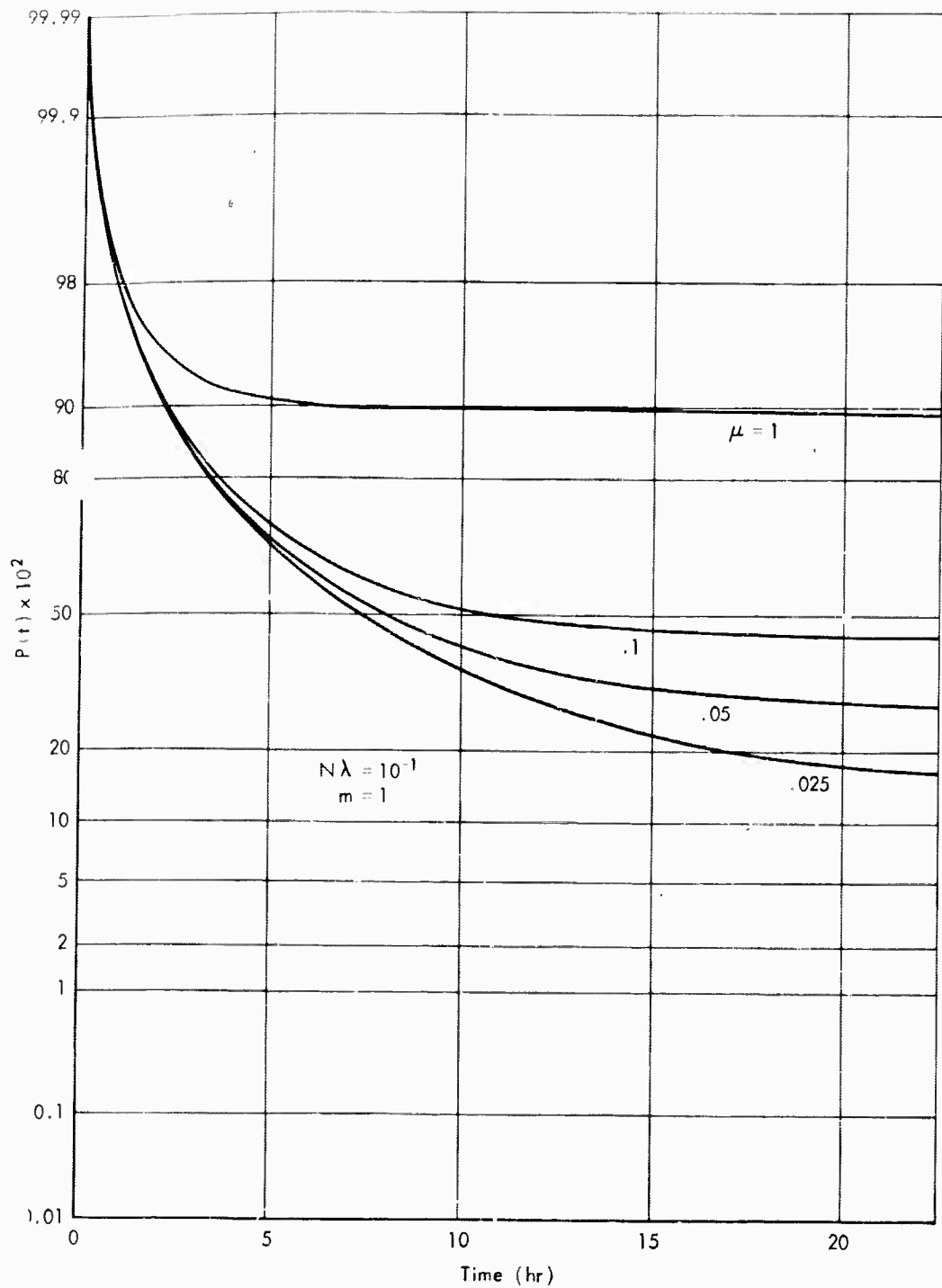


Fig.A-15—Availability of redundant computer
(transient phase, exponential service)

where

$$f^*(s) = \int_0^{\infty} e^{-st} f(t) dt \quad \text{and} \quad g^*(s) = \int_0^{\infty} e^{-st} g(t) dt . \quad (51)$$

Substituting Eq. (51) into (50) and taking the limit yields

$$P_{\infty} = \frac{\int_0^{\infty} t f(t) dt}{\int_0^{\infty} t f(t) dt + \int_0^{\infty} t g(t) dt} = \frac{E[T_f]}{E[T_f] + E[T_s]} . \quad (52)$$

For the non-redundant case of Eq. (30) for which $f(t) = N\lambda \exp(-N\lambda t)$ and $g(t) = \mu \exp(-\mu t)$, we get the result $P_{\infty} = \mu / (N\lambda + \mu)$ as shown by Eq. (30). For the redundant case at hand, $m = 1$, $E[T_f] = 5M/6N\lambda$, and $E[T_s] = 1/\mu$. Hence

$$P_{\infty} = \frac{5M\mu}{5M\mu + 6N\lambda} . \quad (53)$$

Then P_{∞}^M , the asymptotic availability of the entire computer, is shown for various values of $N\lambda$, M , and μ in Figs. A-16 to A-19.

Next, let $m = 2$ (five-fold redundancy). The formulation is exactly the same as $m = 1$, except evaluation of Eq. (40) gives

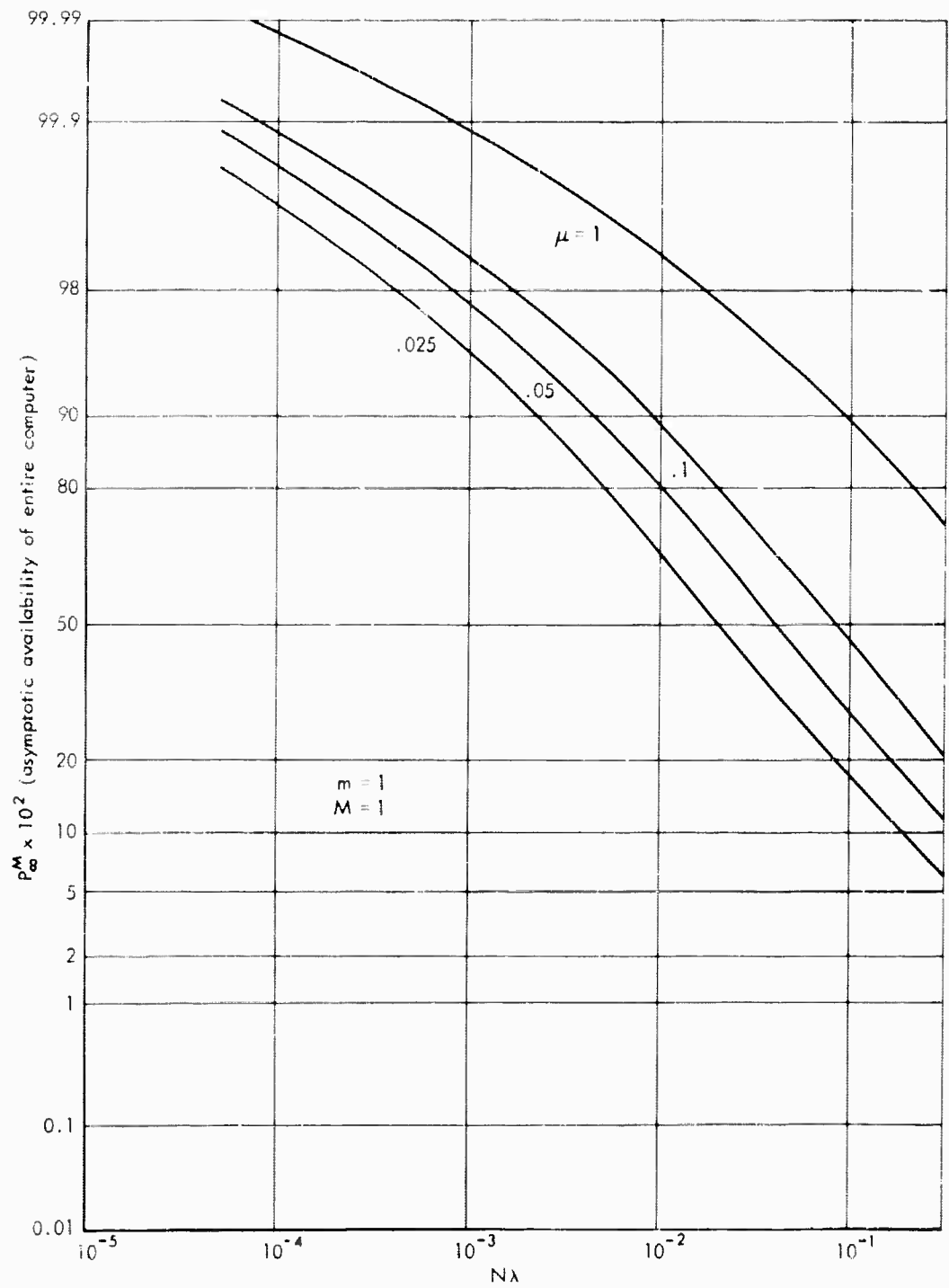


Fig.A-16 — Asymptotic availability of redundant computer
(exponential service)

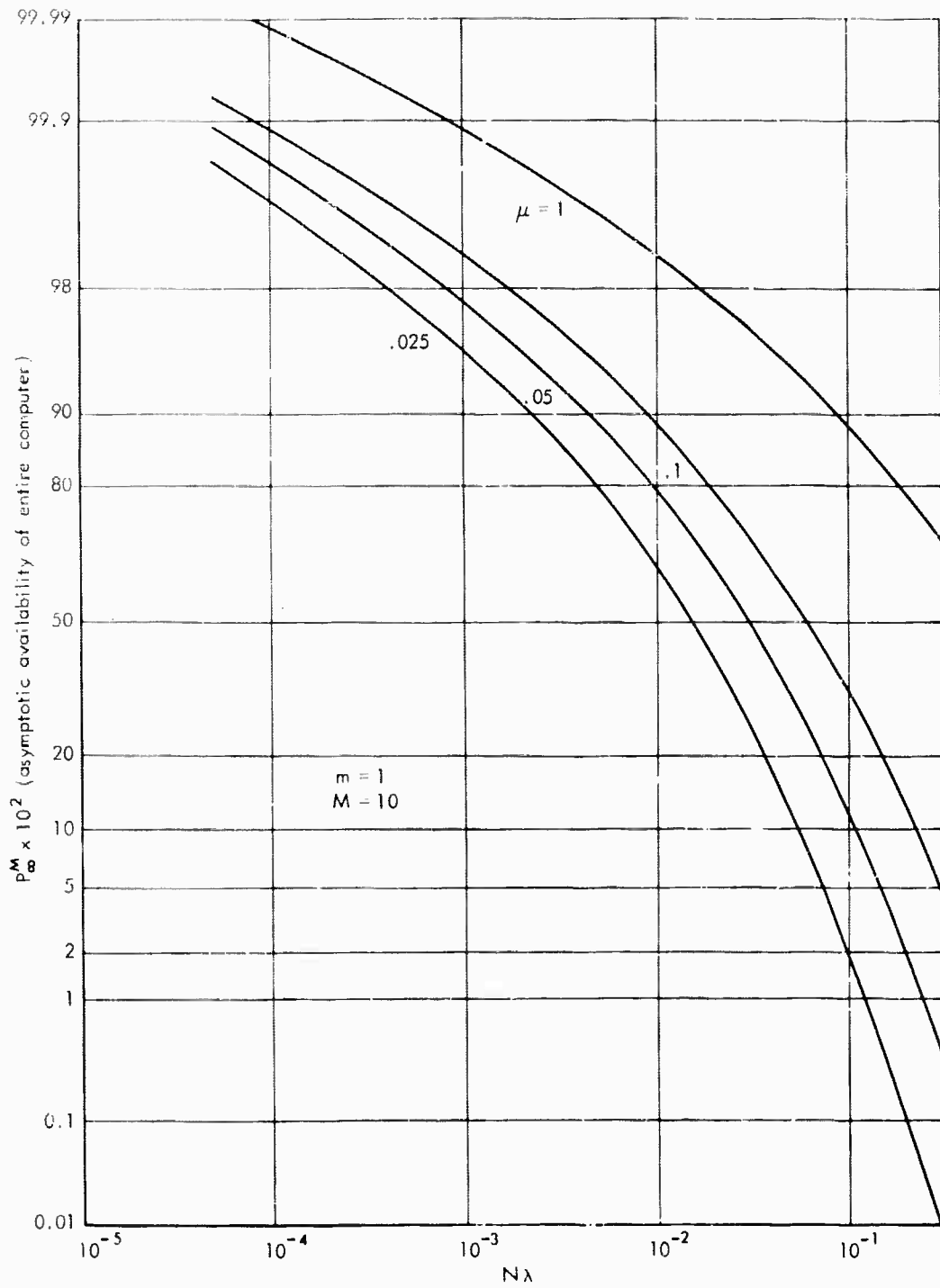


Fig.A-17—Asymptotic availability of redundant computer
(exponential service)

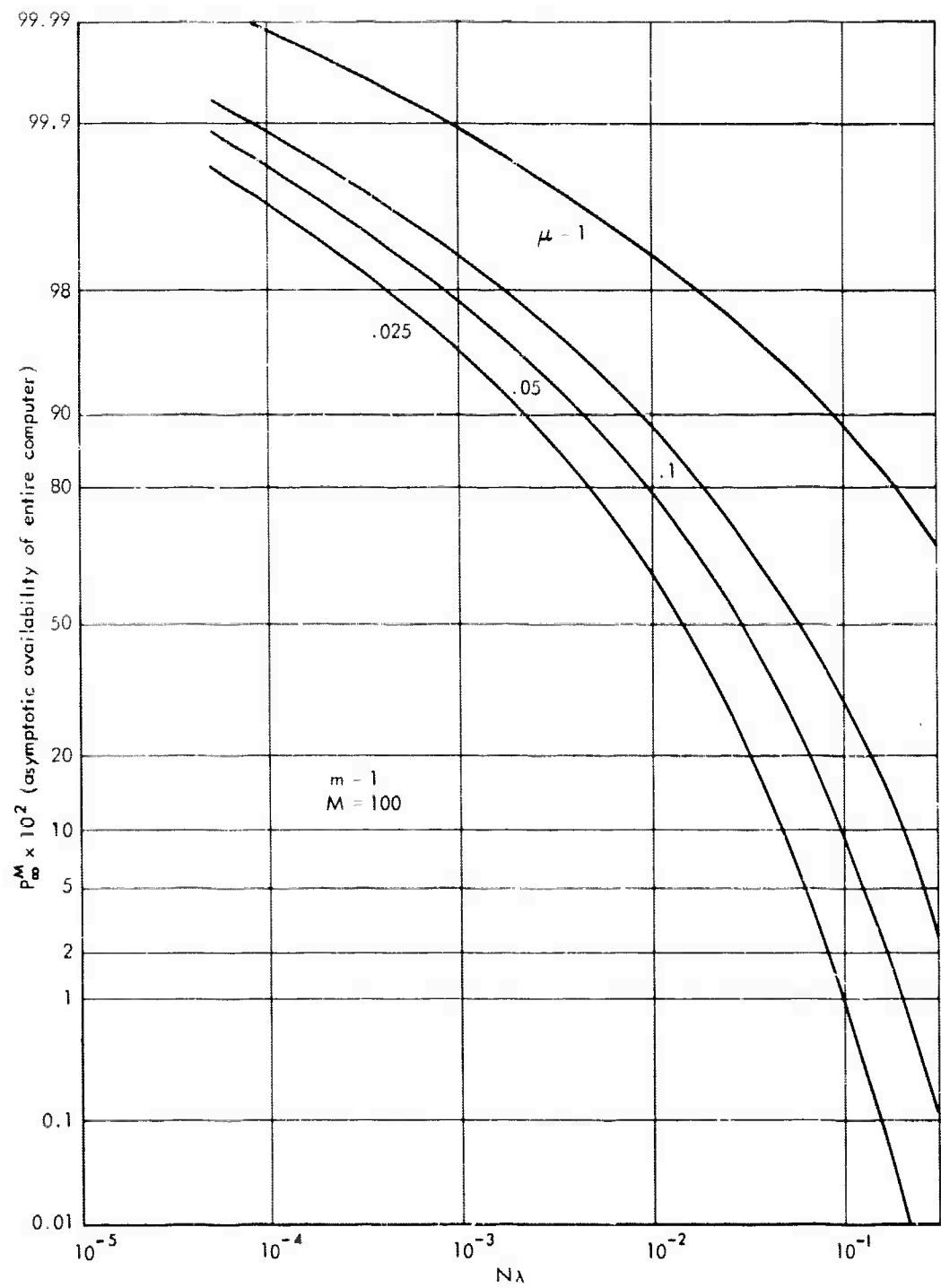


Fig.A-18—Asymptotic availability of redundant computer
(exponential service)

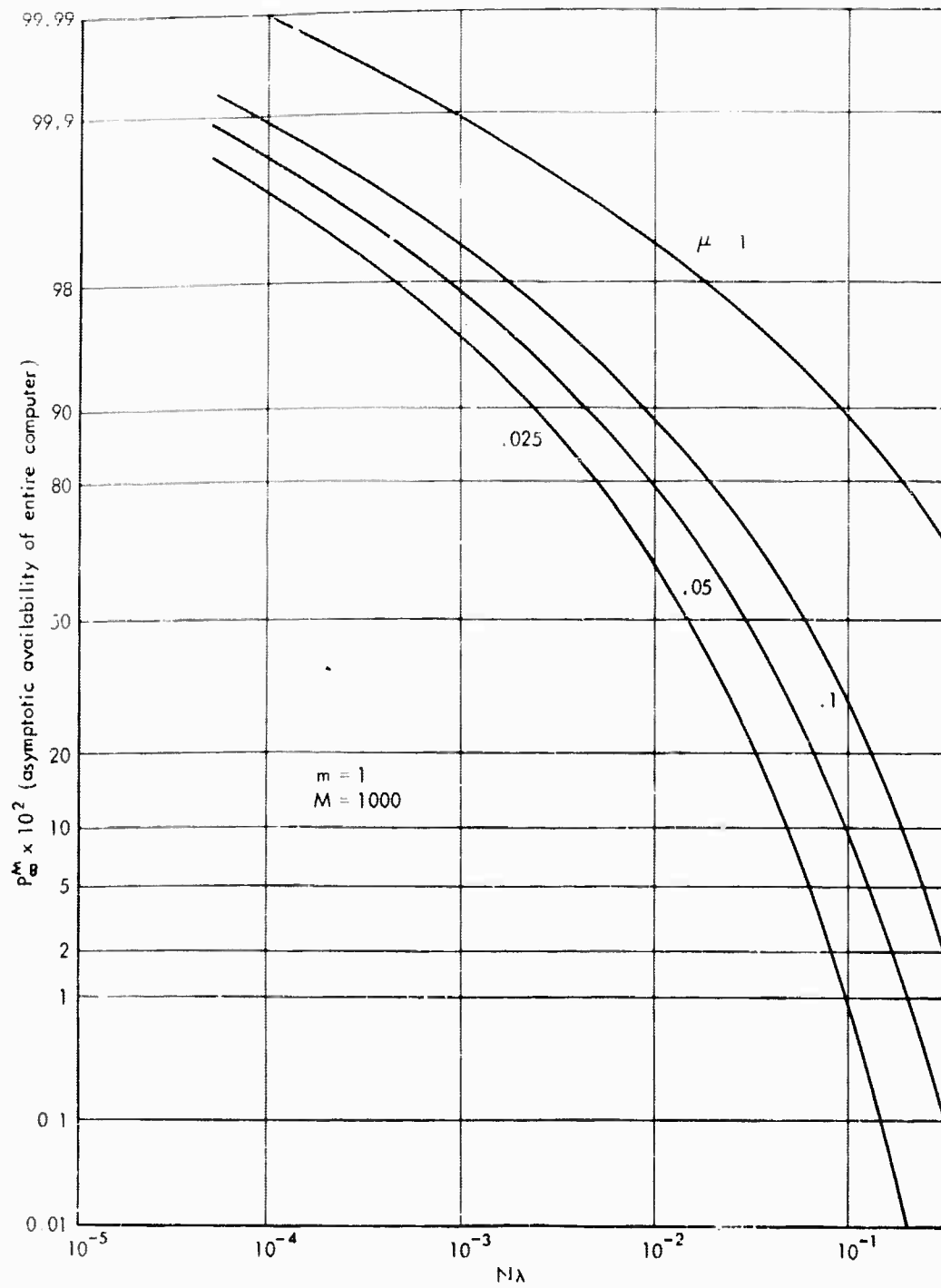


Fig.A-19—Asymptotic availability of redundant computer
(exponential service)

$$F(t) = 1 - 10e^{-3\eta t} + 15e^{-4\eta t} - 6e^{-5\eta t} \quad (54)$$

differentiating

$$f(t) = 30\eta (e^{-3\eta t} - 2e^{-4\eta t} + e^{-5\eta t}) \quad (55)$$

and, as before, $g(t) = \mu \exp(-\mu t)$.

For $m = 2$, the Laplace transform method will, in theory, work, but the effort doesn't appear warranted; only the asymptotic case will be examined. From Eq. (55), we find $E[T_f] = 47/60\eta$. Then, from Eq. (52),

$$P_{\infty} = \frac{47M\mu}{47M\mu + 60N\lambda} \quad (56)$$

Figures A-20 to A-23 show P_{∞}^M for $m = 2$.

Before leaving the subject of transient solutions, presenting a numerical method for computing $P(t)$ is worthwhile, in case no analytical procedure can be found. This method has pitfalls--the user should be warned that the choice of a time step, Δ , which might be required to yield a sufficiently accurate (or even convergent) solution may be exceedingly small for certain choices of $N\lambda$, M , and μ .

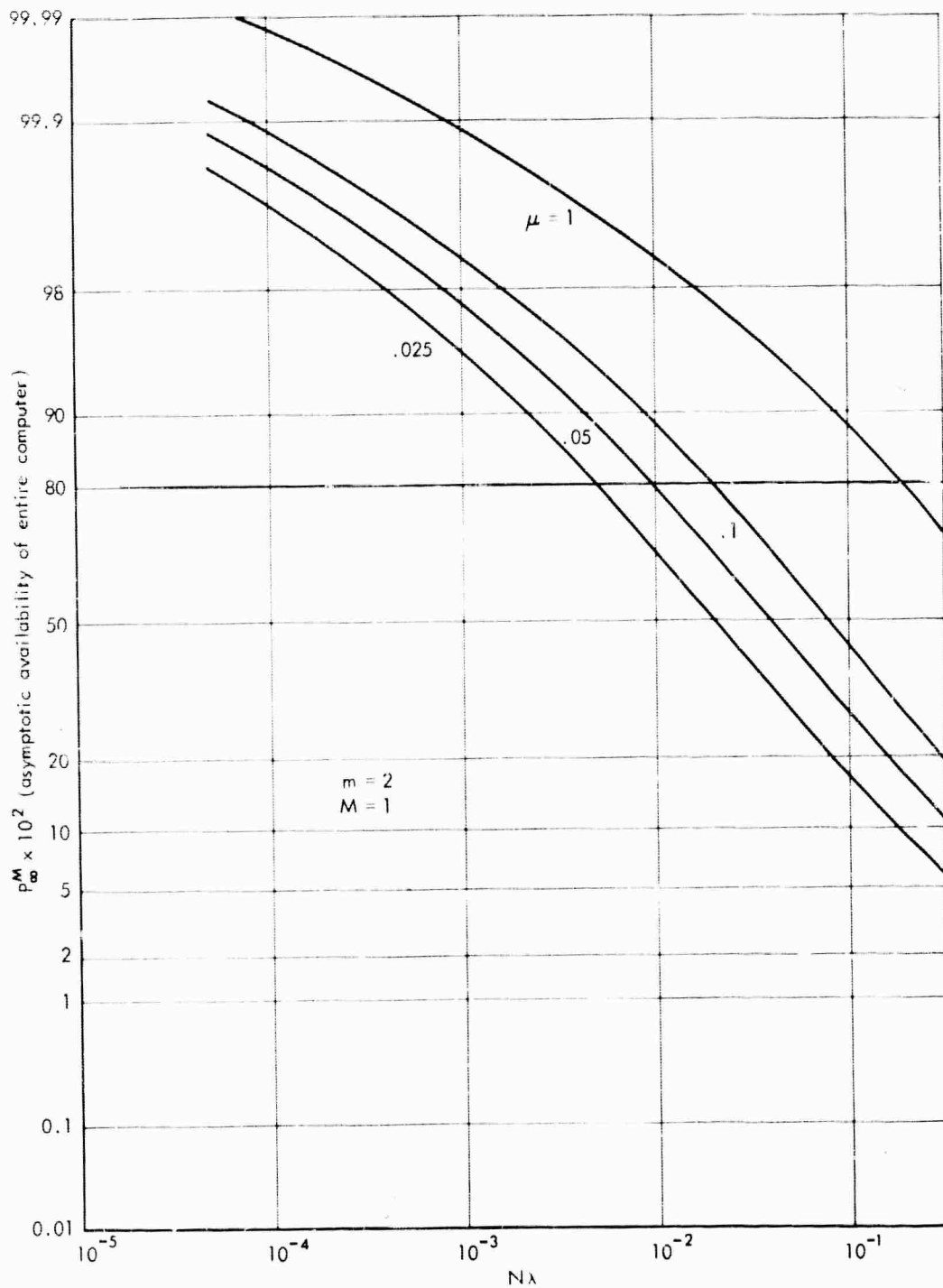


Fig.A-20—Asymptotic availability of redundant computer
(exponential service)

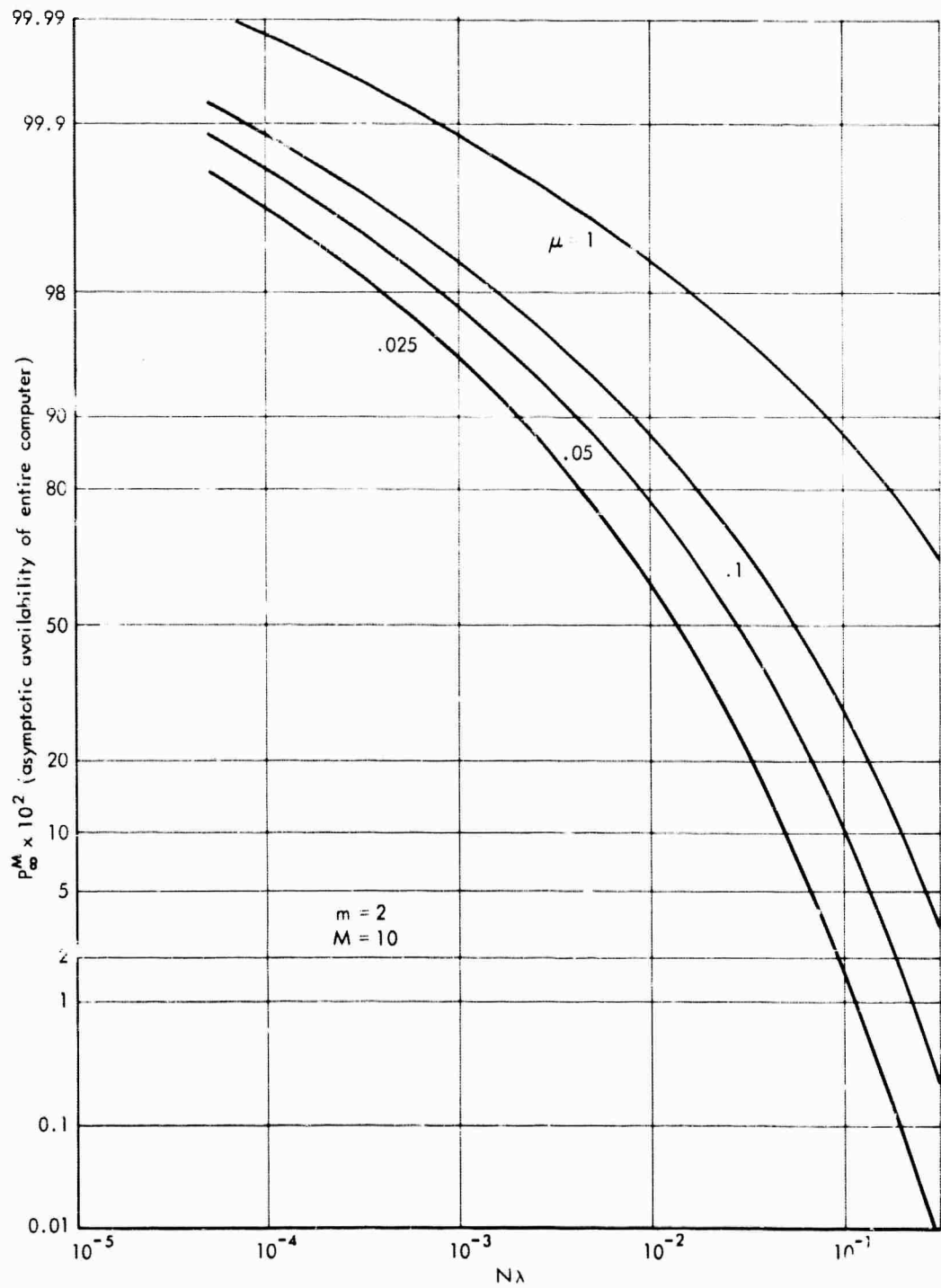


Fig. A-21 — Asymptotic availability of redundant computer
(for sequential service)

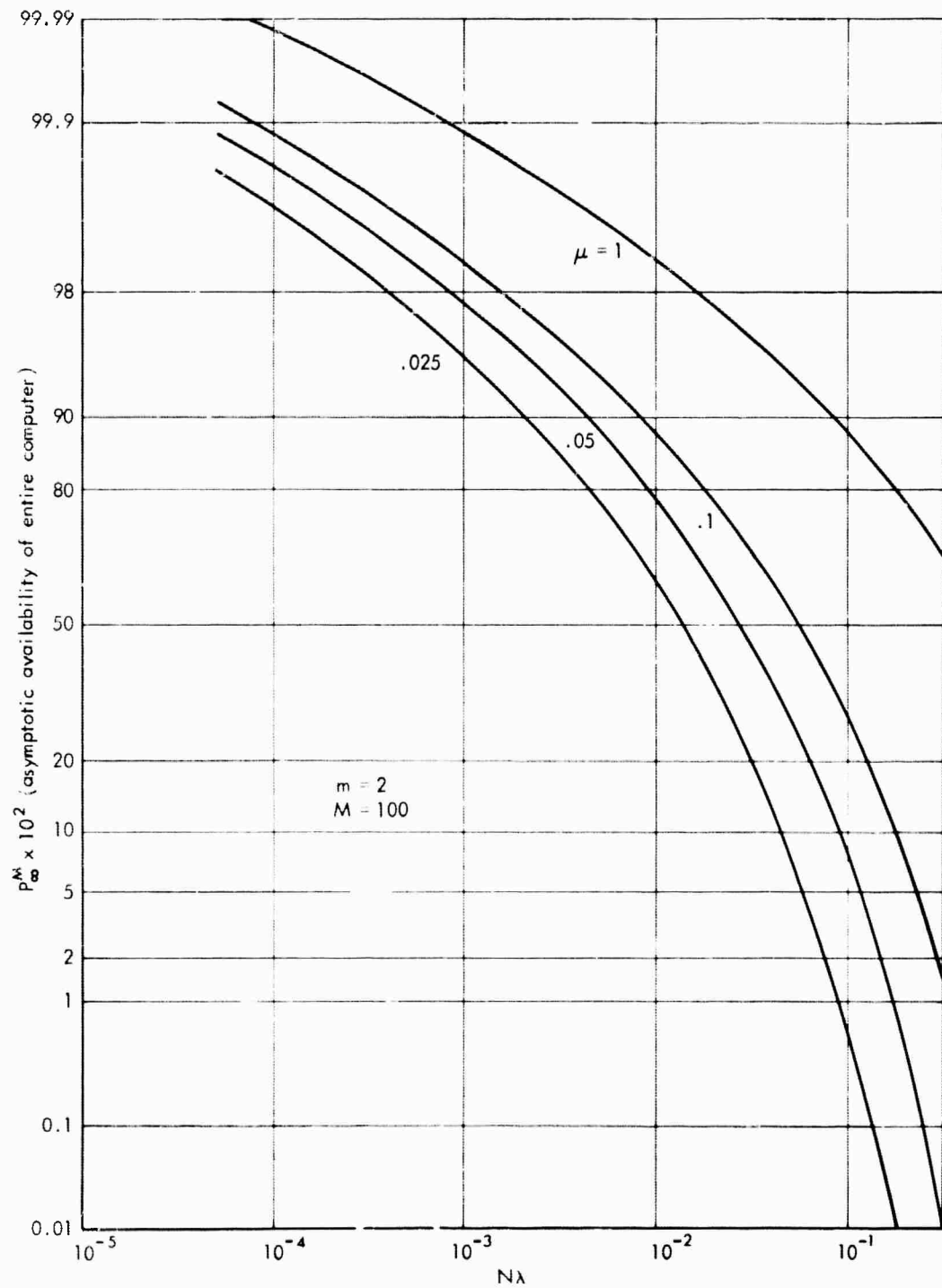


Fig.A-22 — Asymptotic availability of redundant computer
(exponential service)

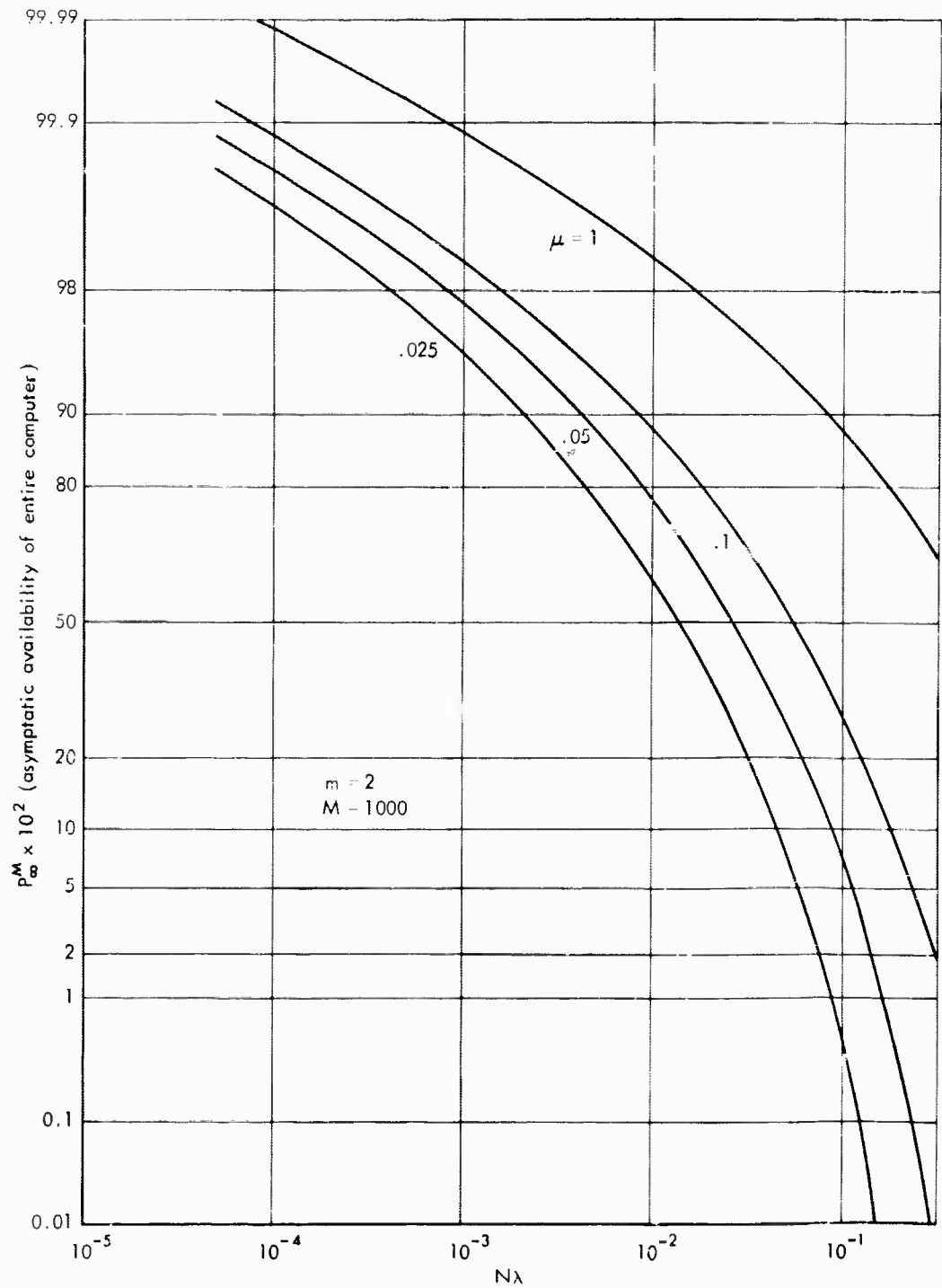


Fig.A-23 — Asymptotic availability of redundant computer
(exponential service)

No guide in making this choice can be given, so a certain amount of trial and error should be expected.[†]

Start again with the Volterra equation (Eq. (37))

$$P(t) = 1 - F(t) + \int_0^t P(u)h(t-u)du \quad (57)$$

where

$$F(t) = \int_0^t f(u)du \quad (58)$$

and

$$h(t-u) = \int_0^{t-u} f(t-u-v)g(v)dv . \quad (59)$$

Define the discrete variables $P_k = P(k\Delta)$, $h_i = h[(k-j)\Delta]$, where $i=k-j$ and $k, j=1,2,\dots$. Assume $F(t)$ has already been evaluated. Then performing the integration in Eq. (57) by the trapezoidal rule gives

[†]For example, in the transient solution of $m=1$, $\eta=.01$, $M=1$, and $\mu=.05$, a choice of $\Delta=0.4$ hr gave a three-significant-figure accuracy. This can be much worse for smaller values of η .

$$\begin{aligned}
 P_k &= 1 - F_k + \frac{\Delta}{2} \sum_{j=0}^{k-1} (P_j h_{k-j} + P_{j+1} h_{k-j-1}) \\
 &= 1 - F_k + \frac{\Delta}{2} P_0 h_k + \Delta \sum_{j=1}^{k-1} P_j h_{k-j} + \frac{\Delta}{2} P_k h_0 . \quad (60)
 \end{aligned}$$

Solving for P_k gives

$$\begin{aligned}
 P_k &= \frac{2 - 2F_k + \Delta P_0 h_k + 2\Delta \sum_{j=1}^{k-1} P_j h_{k-j}}{2 - \Delta h_0} . \\
 j &< k, \quad k=1,2,\dots \quad (61)
 \end{aligned}$$

From Eq. (59), $h_0 = f_0 g_0$, and, again carrying out the integration,

$$\begin{aligned}
 h_{k-j} &= \frac{\Delta}{2} \sum_{i=0}^{k-j-1} (f_{k-j-i} g_i + f_{k-j-i-1} g_{i+1}) \\
 &= \frac{\Delta}{2} f_{k-j} g_0 + \Delta \sum_{i=1}^{k-j-1} f_{k-j-i} g_i + \frac{\Delta}{2} f_0 g_{k-j} . \quad (62)
 \end{aligned}$$

Combining Eqs. (61) and (62) gives

$$P_k = \frac{2 - 2F_k + \Delta P_0 h_k + \Delta^2 \sum_{j=1}^{k-1} \sum_{i=1}^{k-j-1} p_j (f_{k-j} g_0 + 2f_{k-j-i} g_i + f_0 g_{k-j})}{2 - \Delta h_0} \quad (63)$$

P_0 is the initial condition, the probability that the machine is on at $t=0$, and is presumed known.

8. NON-REDUNDANT COMPUTER--EXPONENTIAL SERVICE AND WEIBULL FAILURE DISTRIBUTIONS

Section III-D discussed the possible use of non-constant failure rates--in particular, the difficulties of testing and accepting the hypothesis that a particular component has, indeed, a failure rate which decreases with time. Setting these difficulties aside for now, we show the consequences of assuming such a decreasing failure rate model.

To date, the most successful distribution (in the sense that it fits the experimental data) with non-constant failure rate has been the Weibull distribution; the link, however, between real physical mechanisms and this distribution is still very tenuous, at best.[†] In any event, the Weibull distribution has density function

[†]Cox [1], pp. 109-110, gives a very interesting relationship between the Weibull and exponential distributions--one with perhaps some relationship to reality.

$$f(t) = N\lambda\alpha t^{\alpha-1} e^{-N\lambda t^\alpha} \quad \alpha > 0 \quad (64)$$

which, for $\alpha = 1$, reduces to the exponential distribution. Choosing α so that $0 < \alpha < 1$ is what determines a decreasing failure rate, and its actual value must be determined by experiment. At this point, the ability to obtain transient results without recourse to strictly numerical methods appears very doubtful. Figure A-24 shows a comparison of the transient phase of an exponential failure machine with one having a Weibull failure distribution with $\alpha = 0.5$ by using the numerical method described in Sec. A-7. These are non-redundant serviced machines with $N\lambda = 0.1$, and $\mu = 1$.

The asymptotic solution proceeds quite easily. First, the expectation of T_f is required:

$$E[T_f] = \int_0^\infty t f(t) dt = N\lambda\alpha \int_0^\infty t^\alpha e^{-N\lambda t^\alpha} dt. \quad (65)$$

Let $u = N\lambda t^\alpha$ to get

$$E[T_f] = \left(\frac{1}{N\lambda}\right)^{1/\alpha} \int_0^\infty u^{1/\alpha} e^{-u} du = \left(\frac{1}{N\lambda}\right)^{1/\alpha} \Gamma(1+1/\alpha). \quad (66)$$

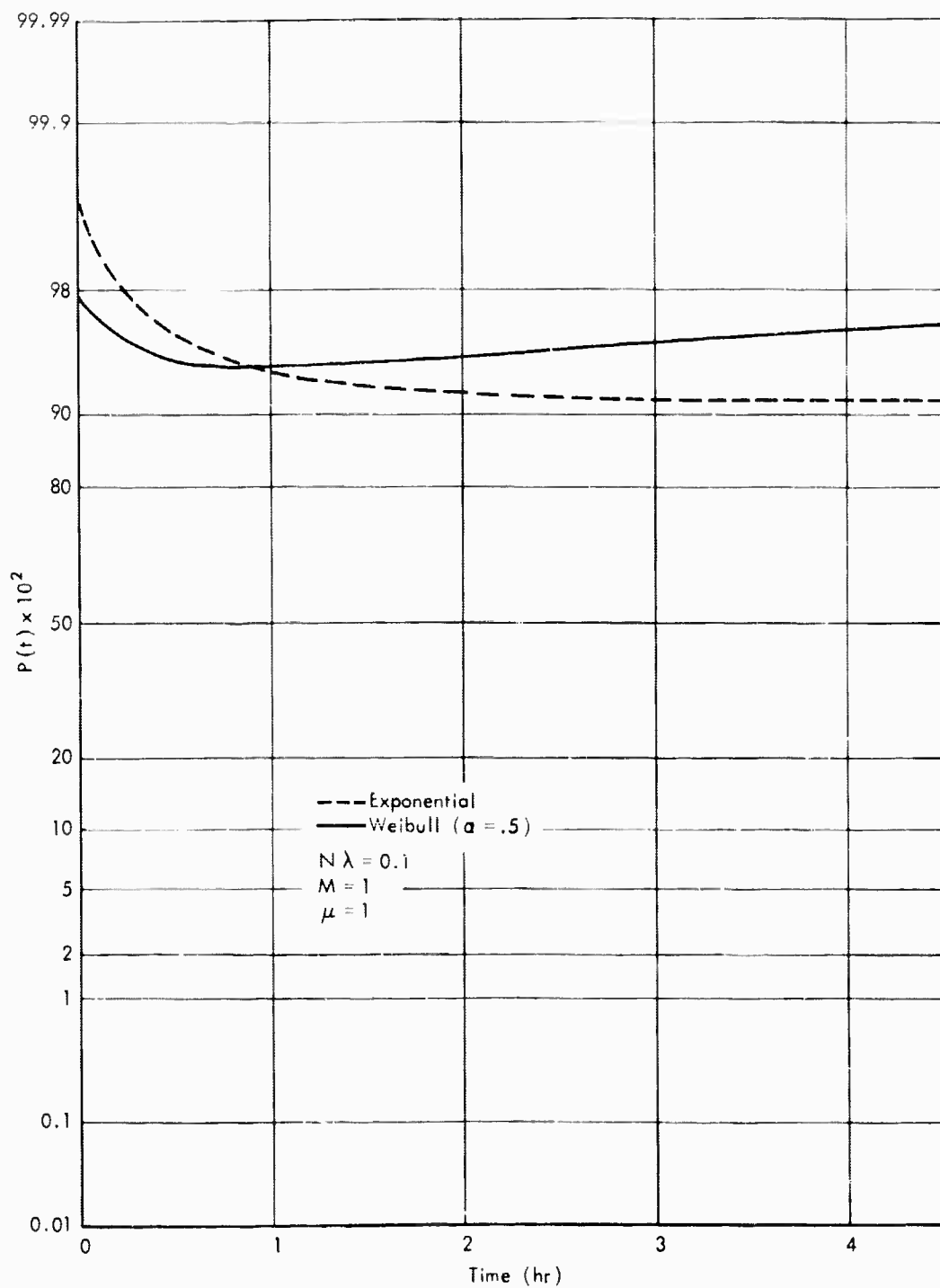


Fig.A-24 — Comparison of transient availability, exponential and Weibull failure distributions (non-redundant computer, exponential service)

Then, assuming as usual that $g(t) = \mu \exp(-\mu t)$,

$$P_{\infty} = \frac{\mu \Gamma(1 + 1/\alpha)}{\mu \Gamma(1 + 1/\alpha) + (N\lambda)^{1/\alpha}}. \quad (67)$$

For the usual values of $N\lambda$ and μ , a comparison between the exponential and Weibull ($\alpha = 0.5$) failure models is shown in Fig. A-25.

9. MULTIPLE COMPUTERS

The next topic is the availability of multiple non-redundant computers. The benefits to be derived from the use of more than one computer, where the extras are treated as spares and profitably used only in the advent of a machine failure, has been discussed in Chap. III; only the analysis appears here.

10. SINGLE COMPUTER--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

This is exactly the case already discussed in Sec. A-6, but the results are repeated here for easy reference. The solution was obtained by formulating the difference equation

$$P_c(t+\Delta) = P_c(t)(1 - N\lambda\Delta) + [1 - P_c(t)]\mu\Delta + o(\Delta), \quad (68)$$

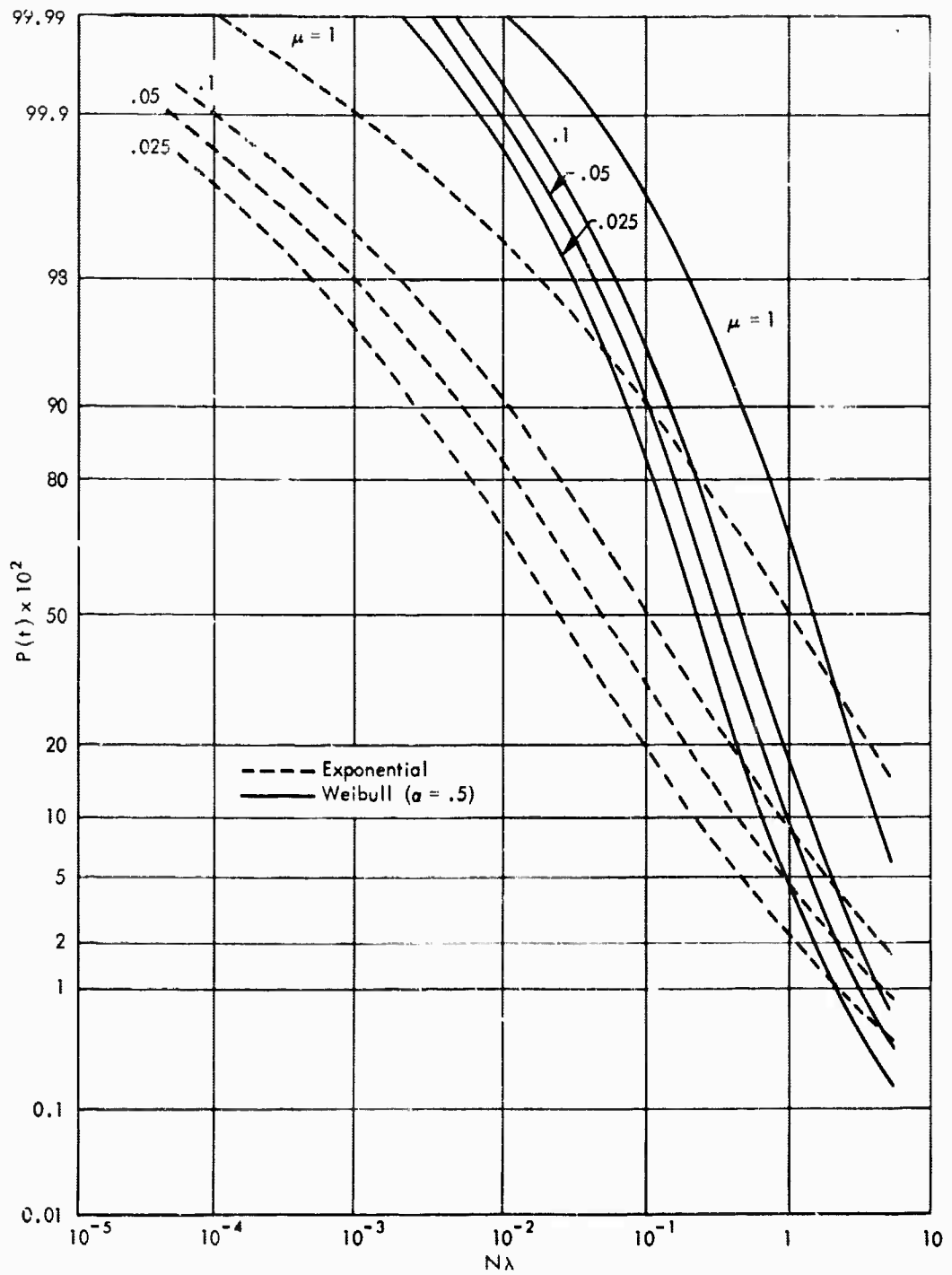


Fig.A-25 — Asymptotic availability of non-redundant computer, exponential and Weibull failure distributions (exponential service)

passing to the limit to get the differential equation

$$P'_c(t) + (N\lambda + \mu)P_c(t) = \mu \quad (69)$$

and solving

$$P_c(t) = \frac{\mu}{N\lambda + \mu} + \frac{N\lambda}{N\lambda + \mu} e^{-(N\lambda + \mu)t} . \quad (70)$$

11. DUPLEX COMPUTERS WITH FLEXIBLE SERVICE--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

When discussing duplex systems, we assume that the system has adequate error checking. Either each machine can check itself or the machines work together to ascertain which is in error.

The simpler of the two service situations has flexible service, so that each machine is independent of the other and the probability that the system has at least one machine operating is easily written as

$$P(t) = 1 - [1 - P_c(t)]^2 \quad (71)$$

where $P_c(t)$, the availability of a single machine system, is given by Eq. (70). In the limit, the asymptotic probability is

$$P_{\infty} = \frac{2\mu}{N\lambda + \mu} - \left(\frac{\mu}{N\lambda + \mu} \right)^2 \quad (72)$$

A graph of P_{∞} is given in Fig. A-26. This should be compared with Fig. A-12 where the asymptotic value of Eq. (70) is shown.

12. DUPLIX COMPUTERS WITH LIMITED SERVICE--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

"Limited service" means that both computers cannot be serviced at the same time. If both fail, one must wait for service until the other is fixed. This introduces a dependence between the two renewal processes which complicates the problem.

Following the derivation in Sec. A-6, and assuming two machines, Computers A and B, $P_a(t)$ is the probability that Computer A is on at time t . $P_b(t)$ has a similar definition. Let α be the probability that Computer A is selected for service if both computers have failed. Then Computer B is selected, given that both have failed, with probability $1-\alpha$.

The difference equations now take the form

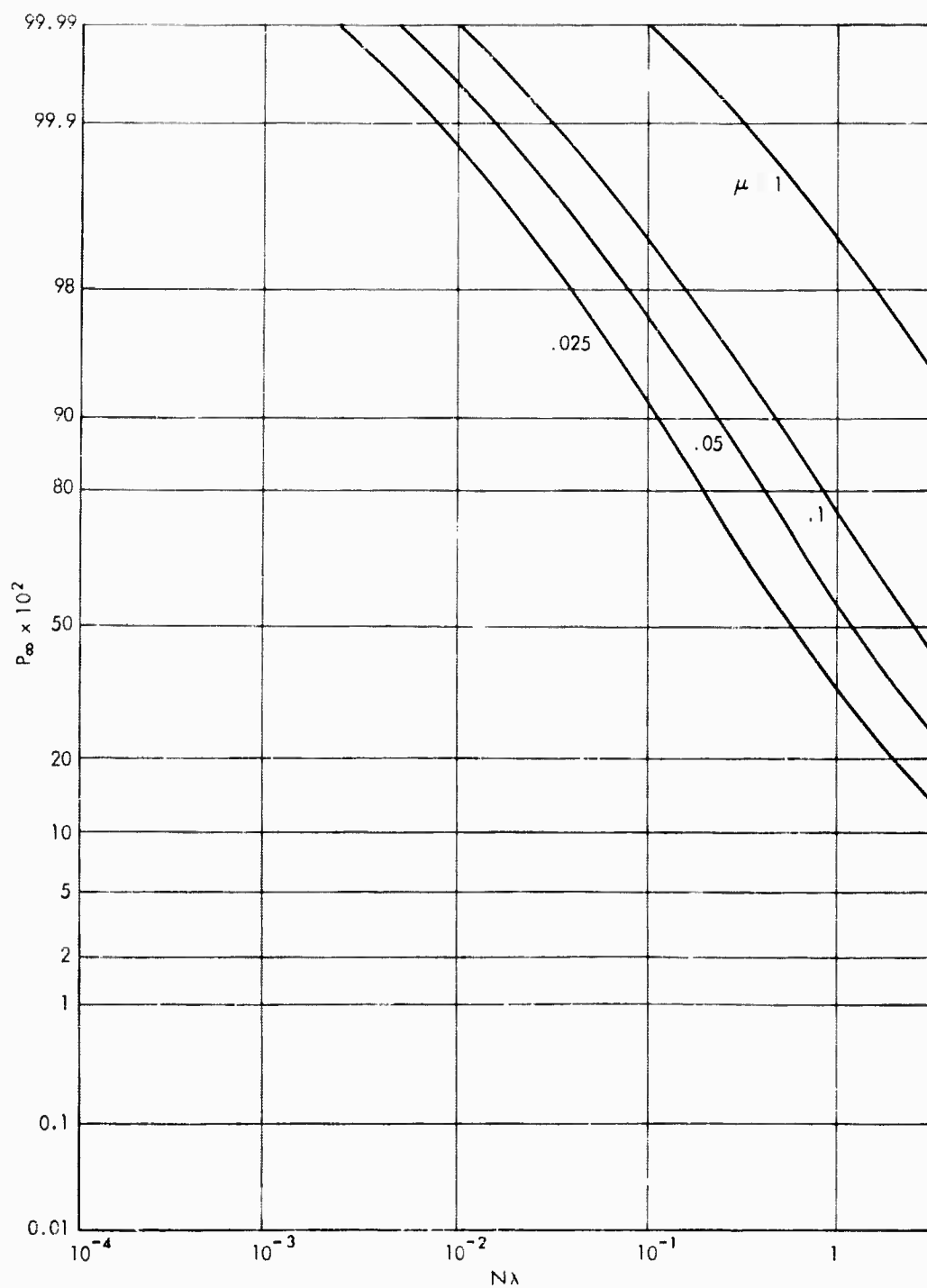


Fig.A-26 — Asymptotic availability of duplex computers (flexible service)

$$P_a(t+\Delta) = P_a(t)(1 - \lambda\Delta) + [1 - P_a(t)]P_b(t)\mu\Delta \\ + [1 - P_a(t)][1 - P_b(t)]\alpha\mu\Delta, \quad (73)$$

and

$$P_b(t+\Delta) = P_b(t)(1 - \lambda\Delta) + [1 - P_b(t)]P_a(t)\mu\Delta \\ + [1 - P_a(t)][1 - P_b(t)](1 - \alpha)\mu\Delta. \quad (74)$$

Rearrange and take the limit to get

$$P_a' + (\lambda + \alpha\mu)P_a - \mu(1 - \alpha)P_b - \mu[\alpha - (1 - \alpha)P_aP_b] = 0, \quad (75)$$

and

$$P_b' + [\lambda + (1 - \alpha)\mu]P_b - \mu\alpha P_a - \mu(1 - \alpha - \alpha P_aP_b) = 0. \quad (76)$$

The problem becomes much more tractable if $\alpha = 1/2$, there being no preferential computer which is serviced first.

Equations (75) and (76) become

$$P_a' + \left(\lambda + \frac{\mu}{2}\right)P_a - \frac{\mu}{2}P_b - \frac{\mu}{2}(1 - P_aP_b) = 0, \quad (77)$$

and

$$P'_b + \left(\lambda + \frac{\mu}{2}\right) P_b - \frac{\mu}{2} P_a - \frac{\mu}{2} (1 - P_a P_b) = 0 . \quad (78)$$

The symmetry between P_a and P_b now allows the distinction between subscripts to be dropped and we have only one differential equation to solve, in, say, the variable P_e :

$$P'_e + \left(\lambda + \frac{\mu}{2}\right) P_e - \frac{\mu}{2} P_e - \frac{\mu}{2} (1 - P_e^2) = 0 ; \quad (79)$$

or, rearranging,

$$P'_e = - \frac{\mu}{2} P_e^2 - \lambda P_e + \frac{\mu}{2} , \quad (80)$$

which is a form of the Ricatti equation. The solution of this equation is possible because the coefficients of the quadratic form are constants and the method may be found in most texts on differential equations [9]. Only the solution is given here. Define

$$P_1 = \frac{-\lambda + \sqrt{\lambda^2 + \mu^2}}{\mu} \quad (81)$$

and

$$P_2 = \frac{-\lambda - \sqrt{\lambda^2 + \mu^2}}{\mu}. \quad (82)$$

Then the general solution, $P_e(t)$, is

$$\frac{P_e(t) - P_1}{P_e(t) - P_2} = K \exp \left[-\frac{\mu}{2}(P_1 - P_2)t \right] \quad (83)$$

and, if $P_e(0) = 1$,

$$P_e(t) = \frac{P_1(1 - P_2) - P_2(1 - P_1) \exp \left[-\frac{\mu}{2}(P_1 - P_2)t \right]}{(1 - P_2) - (1 - P_1) \exp \left[-\frac{\mu}{2}(P_1 - P_2)t \right]}. \quad (84)$$

$P_e(t)$ is the probability that either of the two computers is on at time t . Hence, the probability that the duplex system is operative at t is again

$$P(t) = 1 - [1 - P_e(t)]^2. \quad (85)$$

In the limit, we get $P_\infty = 2P_1 - P_1^2$ as the asymptotic probability. This probability is shown in Fig. A-27, and should be compared to Figs. A-12 and A-26.

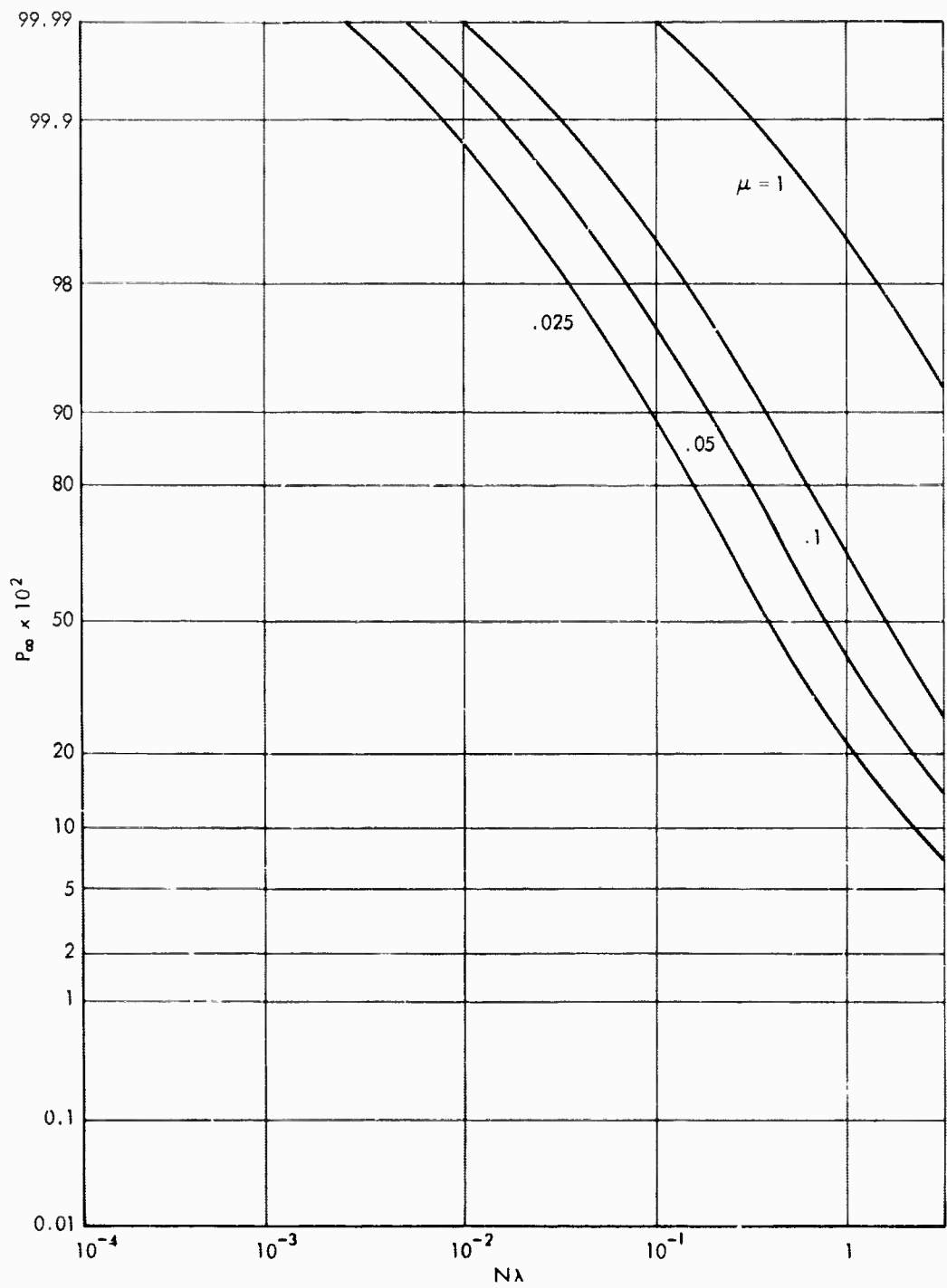


Fig.A-27—Asymptotic availability of duplex computers
(limited service)

13. TRIPLEX COMPUTERS WITH FLEXIBLE SERVICE--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

A system made more reliable by adding a third computer can be operated in two ways. First, adequate error checking is assumed, and therefore the system fails only when all three machines fail. In this case, $P(t) = 1 - [1 - P_c(t)]^3$, where $P_c(t)$ is given by Eq. (70).

If the three machines do not have enough self-checking, their outputs should be majority voted. With this method we have

$$P(t) = \sum_{k=2}^3 \binom{3}{k} P_c^k(t) [1 - P_c(t)]^{3-k} = 3P_c^2(t) - 2P_c^3(t) .$$

(86)

14. THE MULTI-PROCESSOR SYSTEM--EXPONENTIAL SERVICE AND FAILURE DISTRIBUTIONS

Finally, the reliability of the so-called multi-processor must be discussed. Again, the theory underlying the multi-processor structure is given in Sec. IV-2 and only the probabilistic analysis will be done here. The system is shown schematically in Fig. A-28, and assumes the simplification that all the switching circuits are perfectly reliable. If this is not the case, its effect can probably be absorbed into the failure properties of

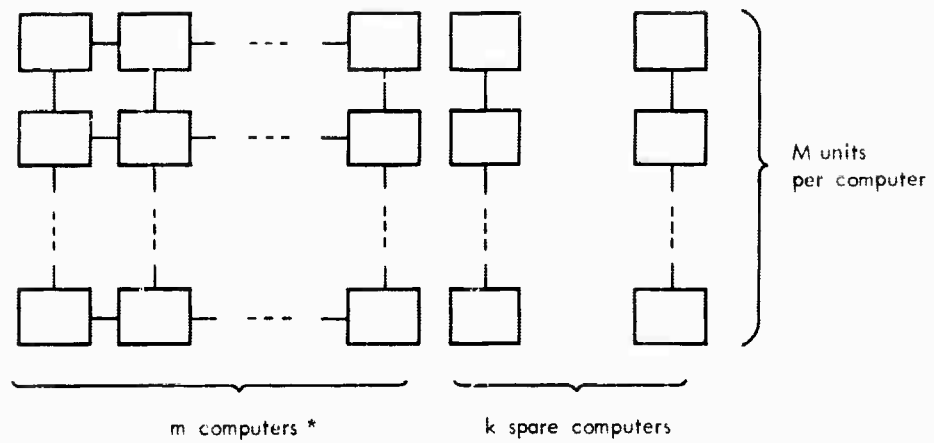


Fig.A-28—The Multi-processor

the individual units. Furthermore, all the individual units are assumed to be identical. Clearly, a memory is not the same as an input/output channel, but for the purposes of this probabilistic analysis it is assumed to be. The following definitions are needed:

M = Number of units needed to make up a single computer,

m = Number of computers need to make up the processor,

N = The total number of parts in the $M \times m$ system,

λ = The failure rate of a single (average) part,

μ = Service rate for a single unit (flexible service),

k = Number of spare units of each type.

Let $P_u(t)$ be the probability that a unit is on at t .

This can be found by replacing N , the number of parts in one unit in Eq. (30), by N/Mm , the number of parts per unit in the multi-processor case, and we get

$$P_u(t) = \frac{Mm\mu}{N\lambda + Mm\mu} + \frac{N\lambda}{N\lambda + Mm\mu} \exp \left[- \left(\frac{N\lambda}{Mm} + \mu \right) t \right] \quad (87)$$

There are enough units of type i only if at least m out of $m+k$ are operating. The probability of this event is

$$\begin{aligned} P_i(t) &= \sum_{j=m}^{m+k} b \left[m+k; j, P_u(t) \right] \\ &= 1 - \sum_{j=0}^{m-1} \binom{m+k}{j} P_u^j(t) [1 - P_u(t)]^{m+k-j} . \quad (88) \end{aligned}$$

Finally, the multi-processor system is on only if all the unit types have at least m of $m+k$ on, and since all the unit types are taken to be identical and independent

$$P(t) = \prod_{i=1}^M P_i(t) = [P_i(t)]^M. \quad (89)$$

Let $P_\infty = \lim_{t \rightarrow \infty} P_i(t)$; i.e., the value of $P_i(t)$ when $P_u = Mm\mu/(N\lambda + Mm\mu)$. Then P_∞^M , the asymptotic availability for the entire processor, is shown for various values of $N\lambda$, M , m , μ , and k in Figs. A-29 to A-40.

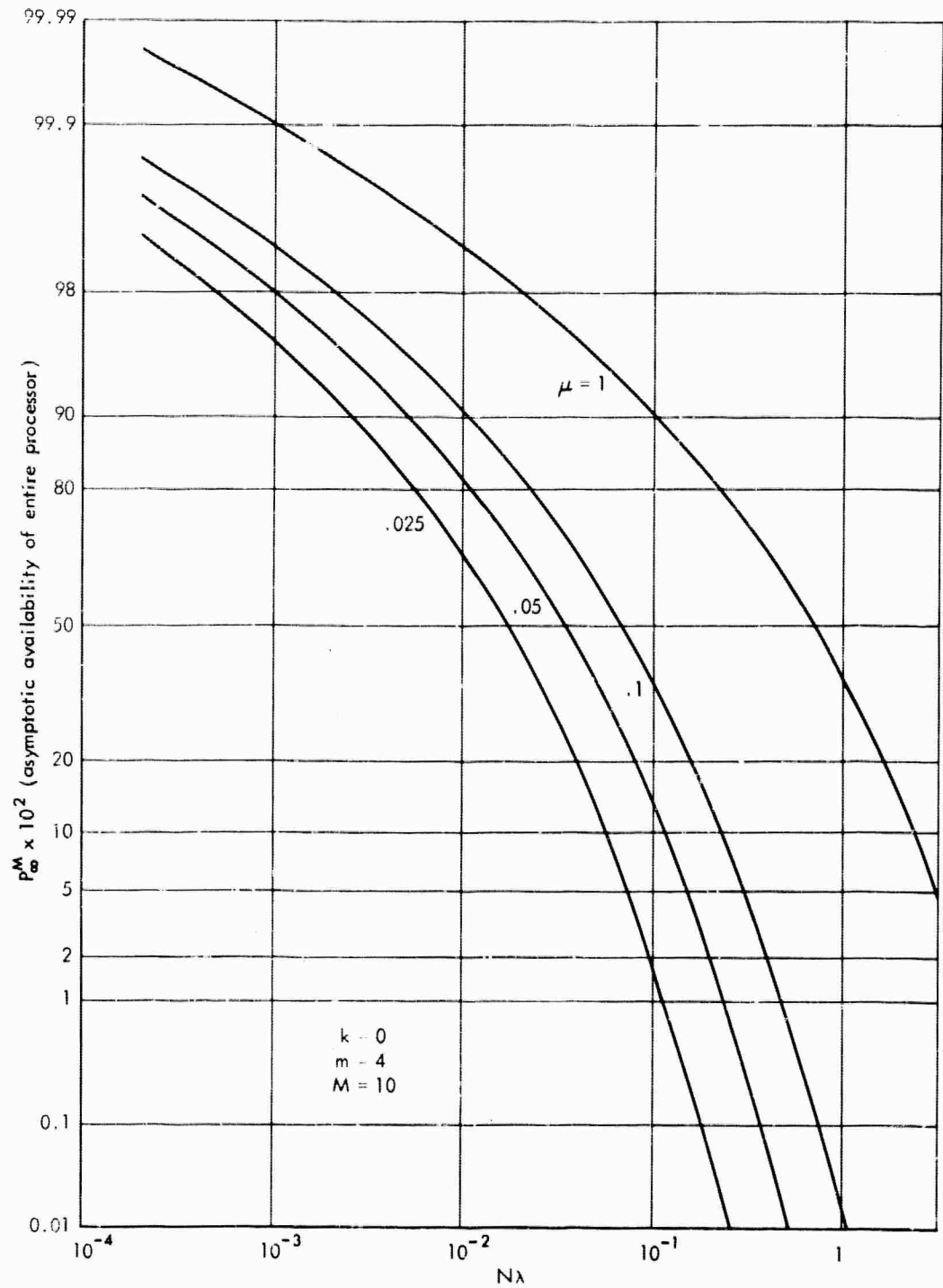


Fig.A-29 — Asymptotic availability of multi-processor
(exponential service)

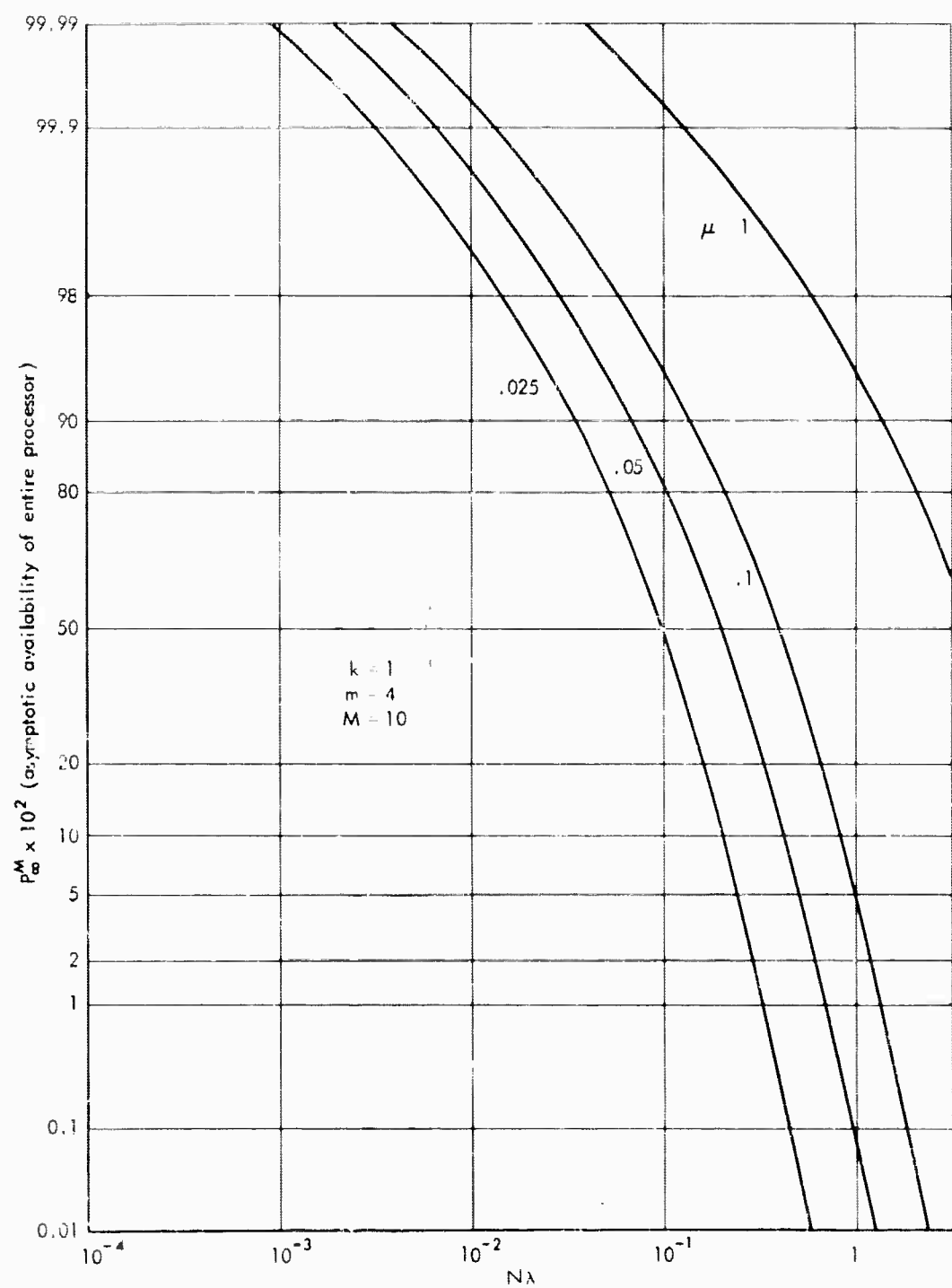


Fig.A-30—Asymptotic availability of multi-processor
(exponential service)

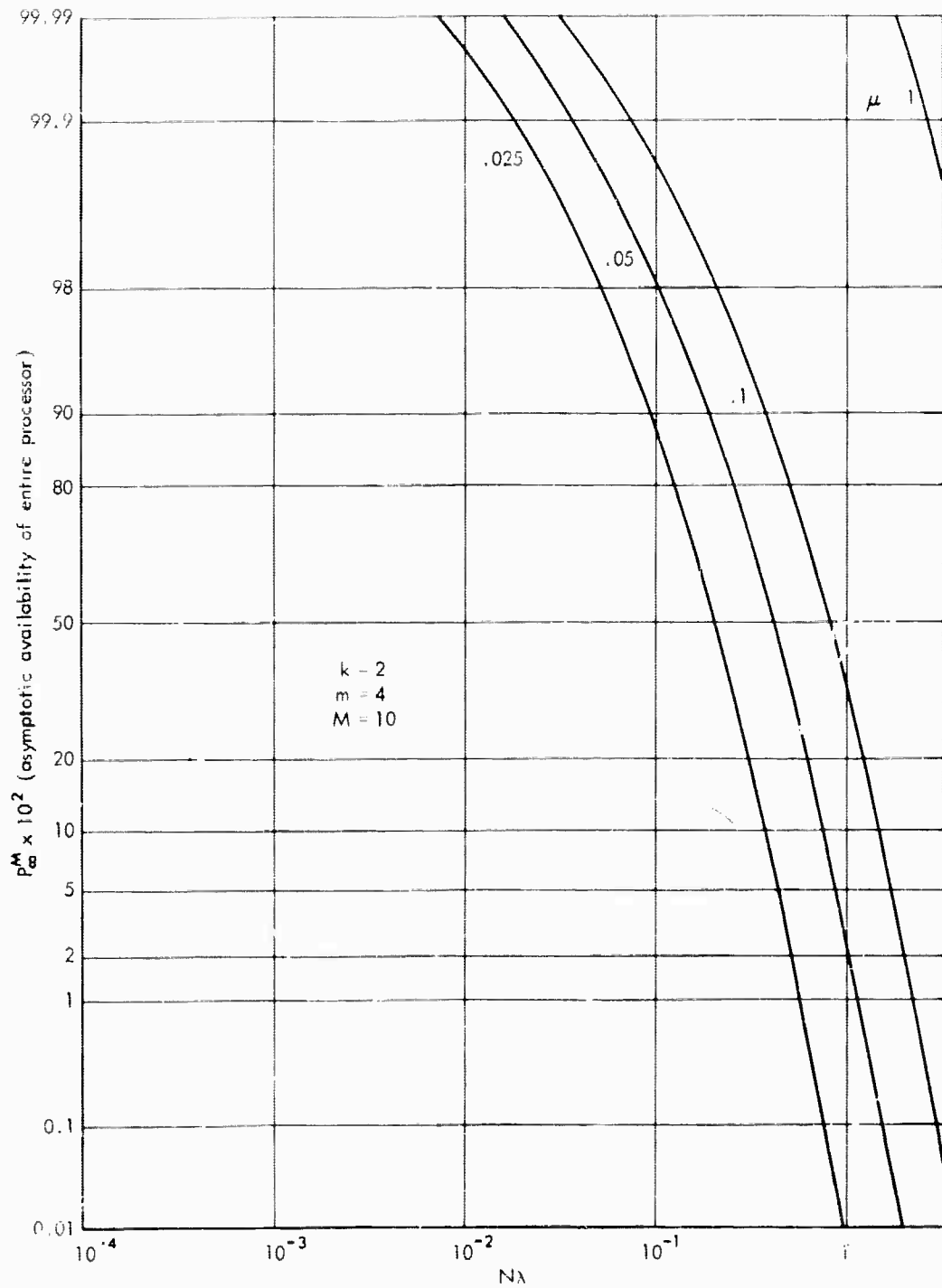


Fig.A-31 — Asymptotic availability of multi-processor
(exponential service)

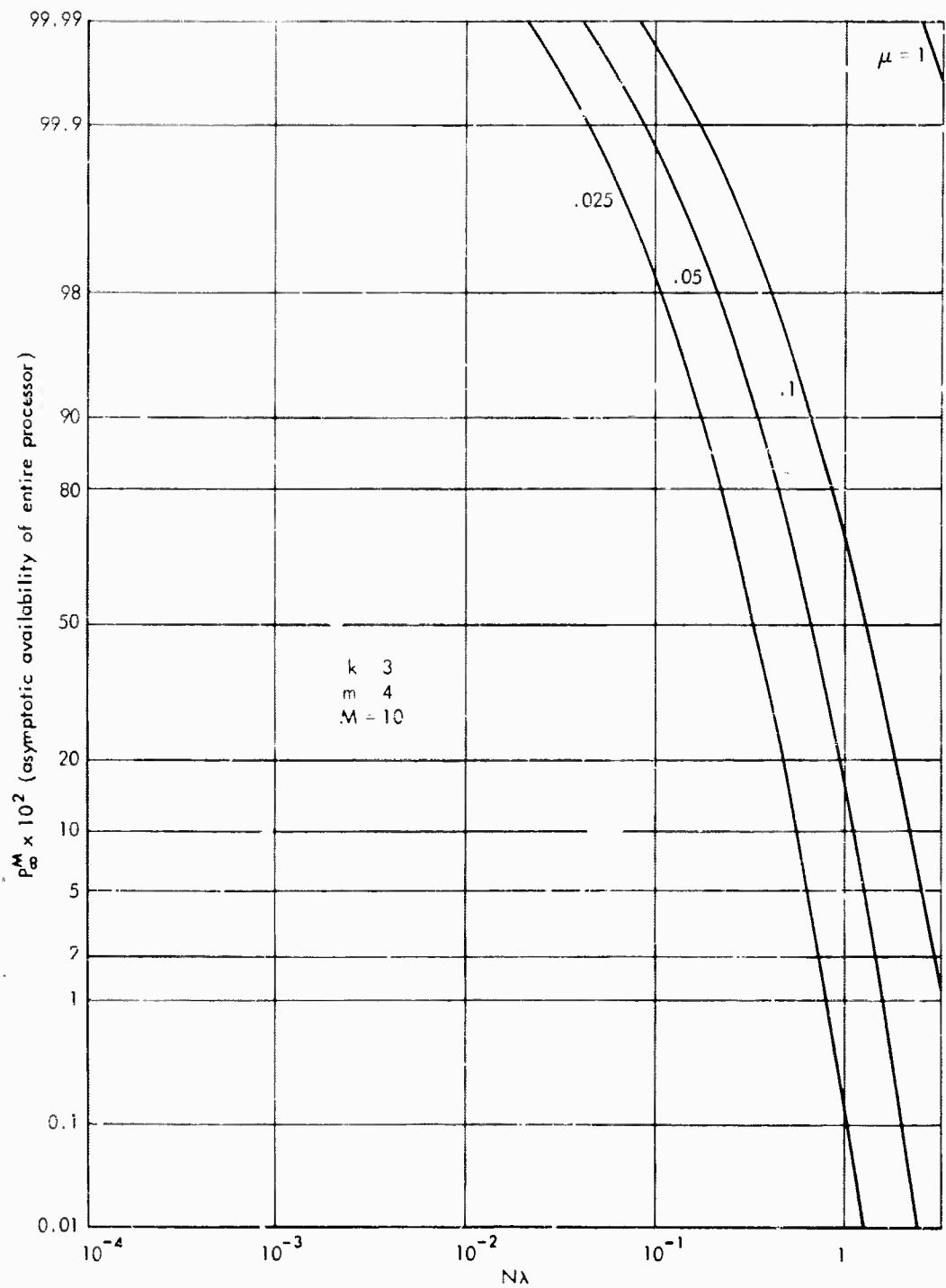


Fig.A-32 — Asymptotic availability of multi-processor (exponential service)

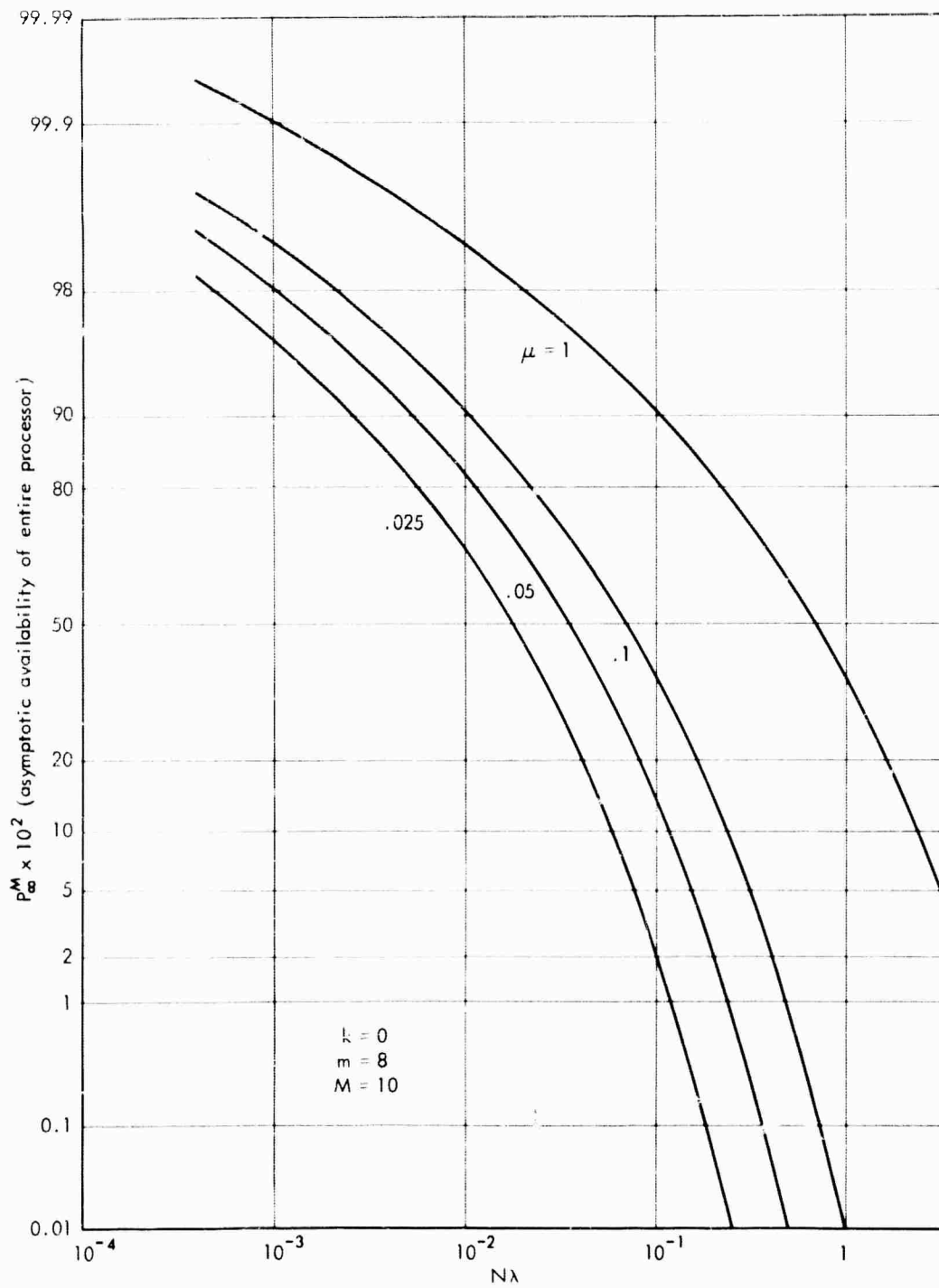


Fig.A-33—Asymptotic availability of multi-processor
(exponential service)

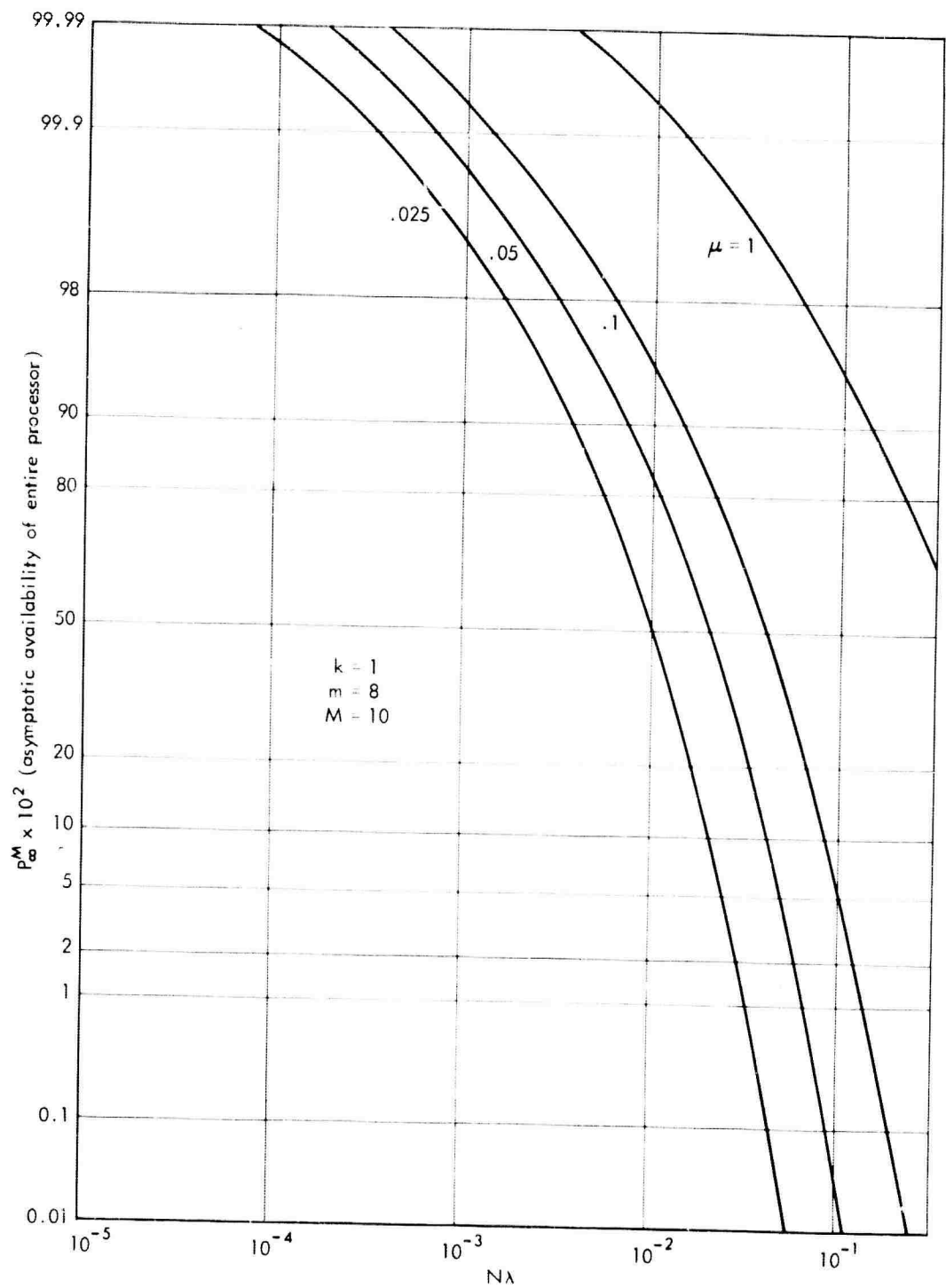


Fig.A-34 — Asymptotic availability of multi-processor (exponential service)

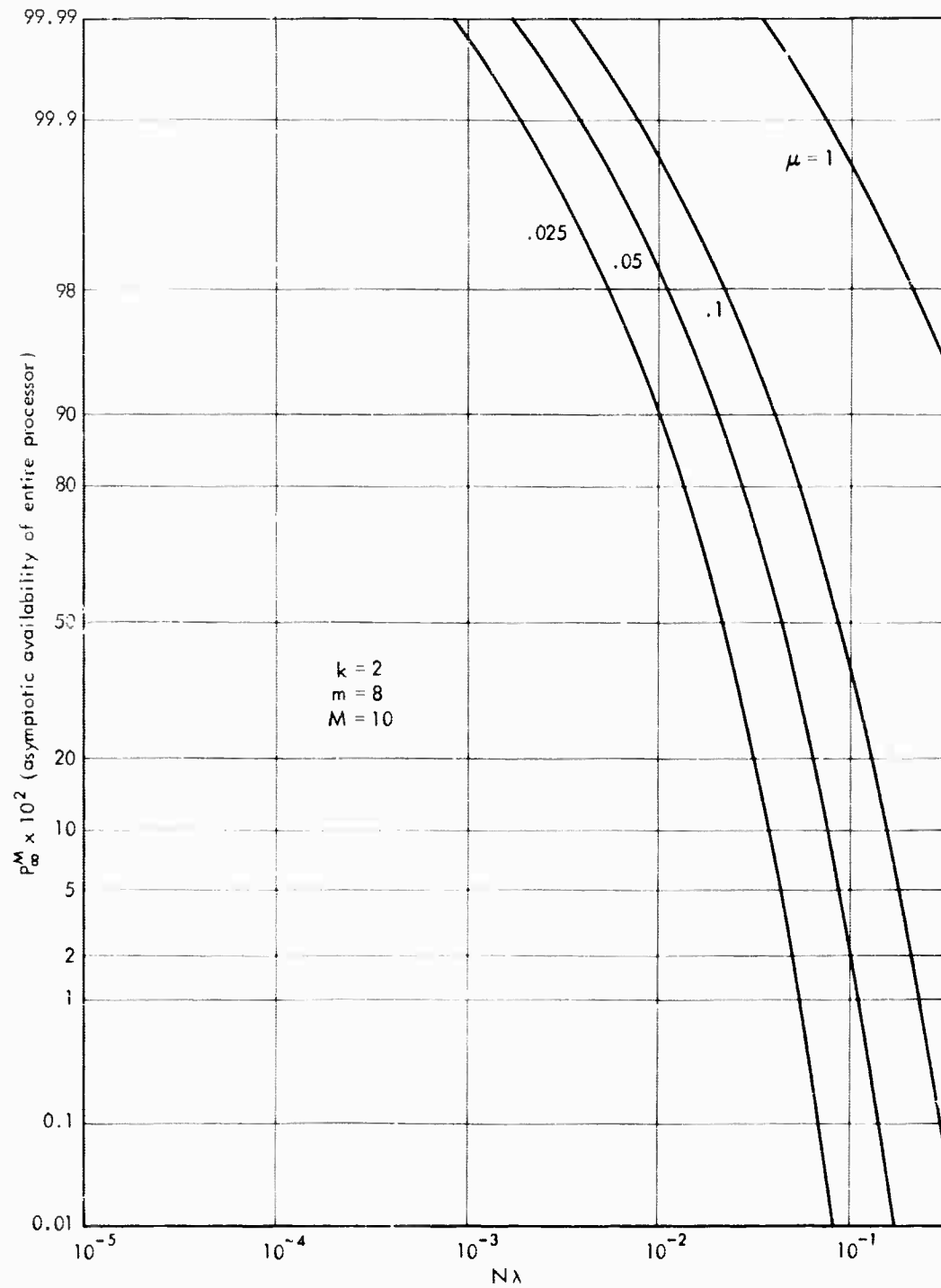


Fig.A-35—Asymptotic availability of multi-processor
(exponential service)

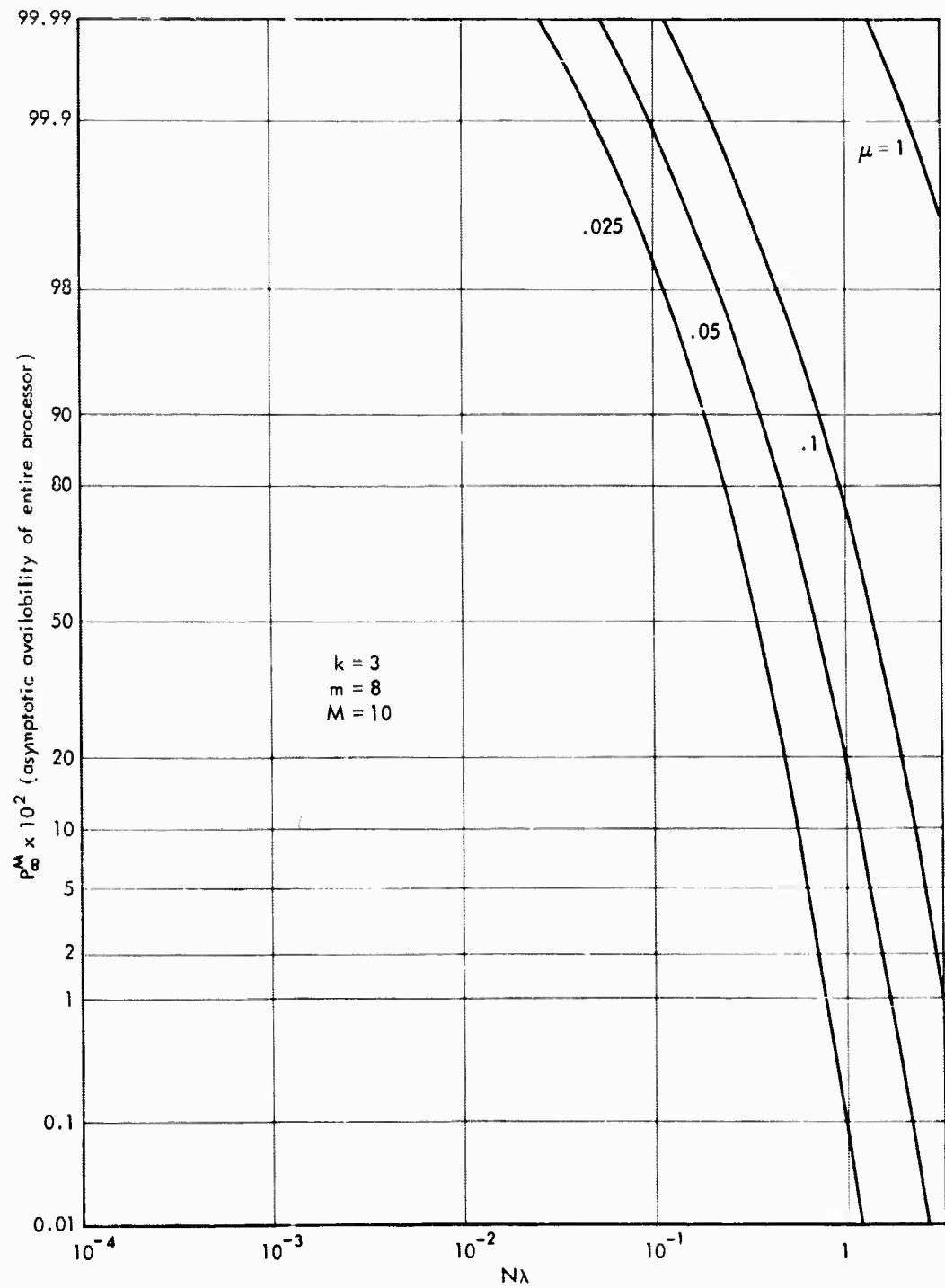


Fig.A-36—Asymptotic availability of multi-processor
(exponential service)

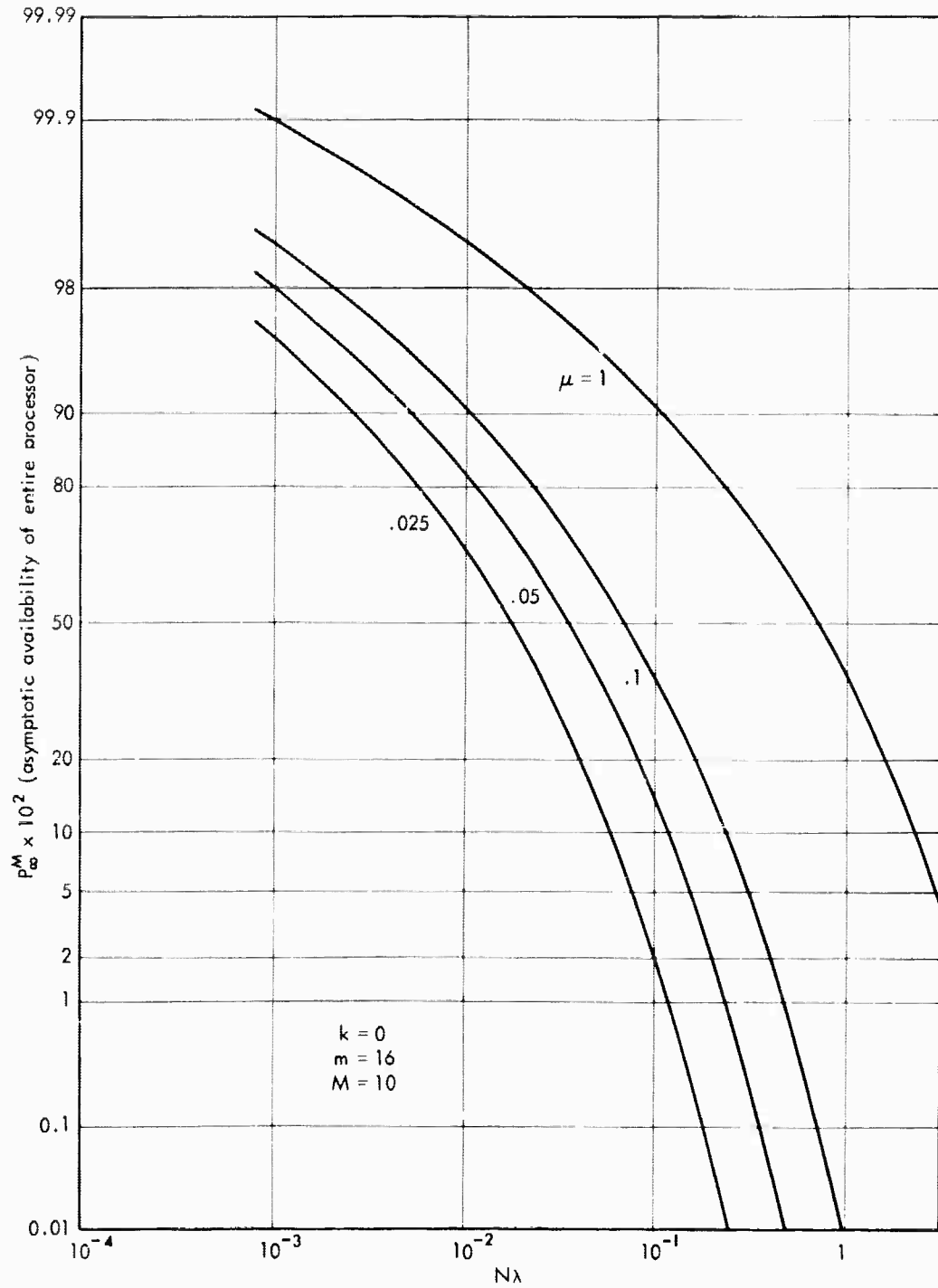


Fig.A-37—Asymptotic availability of multi-processor
(exponential service)

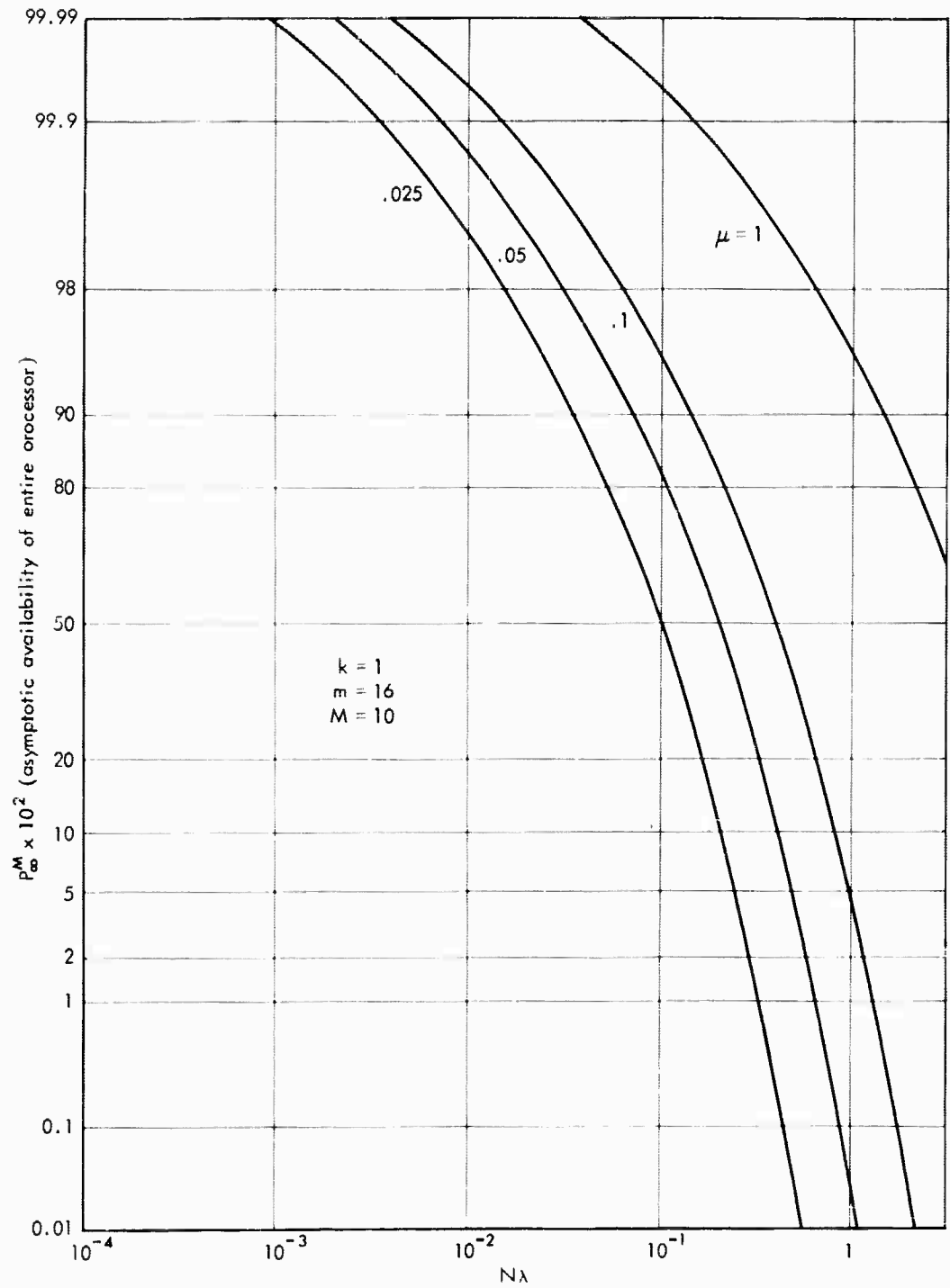


Fig.A-38 — Asymptotic availability of multi-processor
(exponential service)

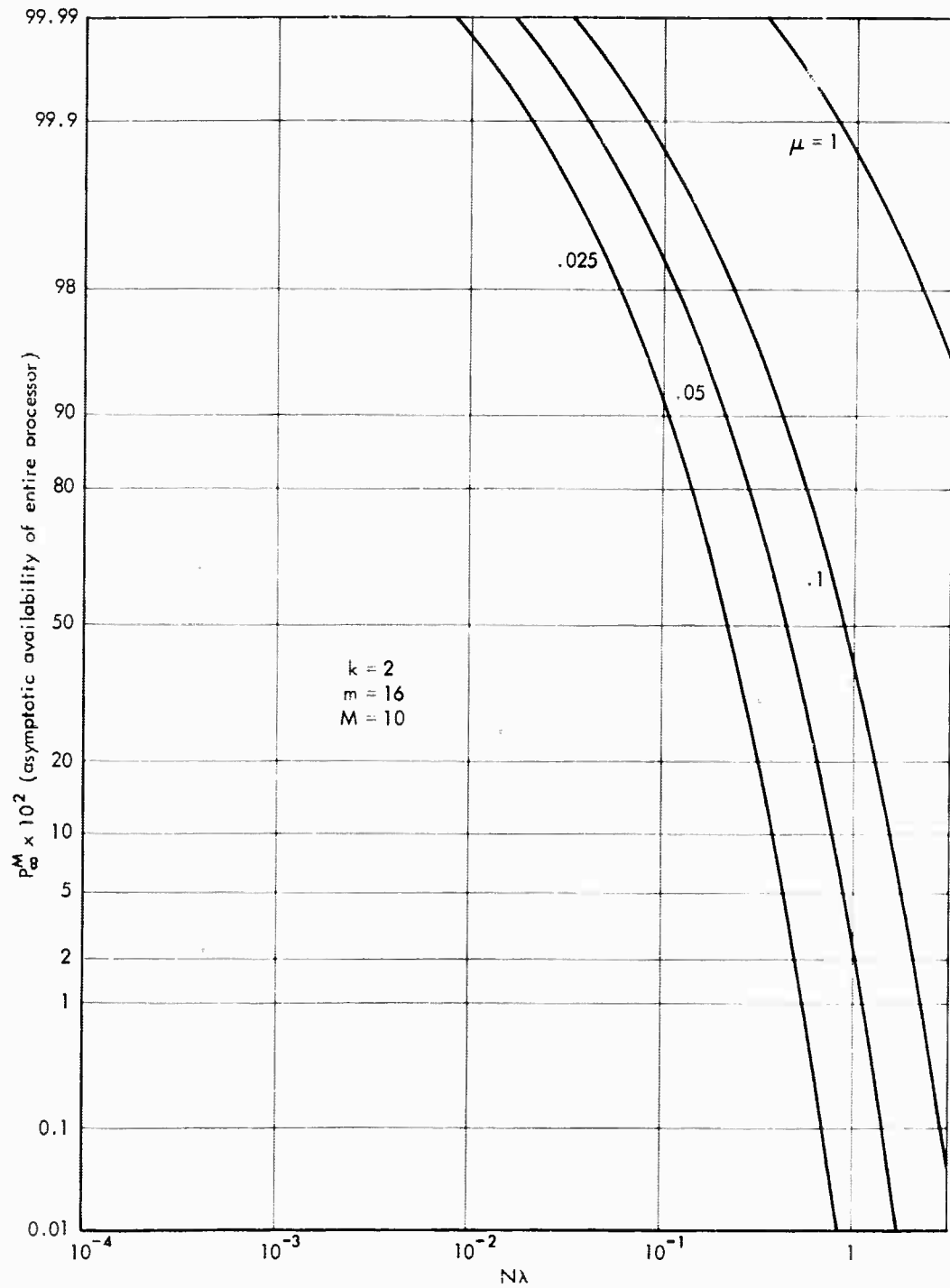


Fig.A-39—Asymptotic availability of multi-processor (exponential service)

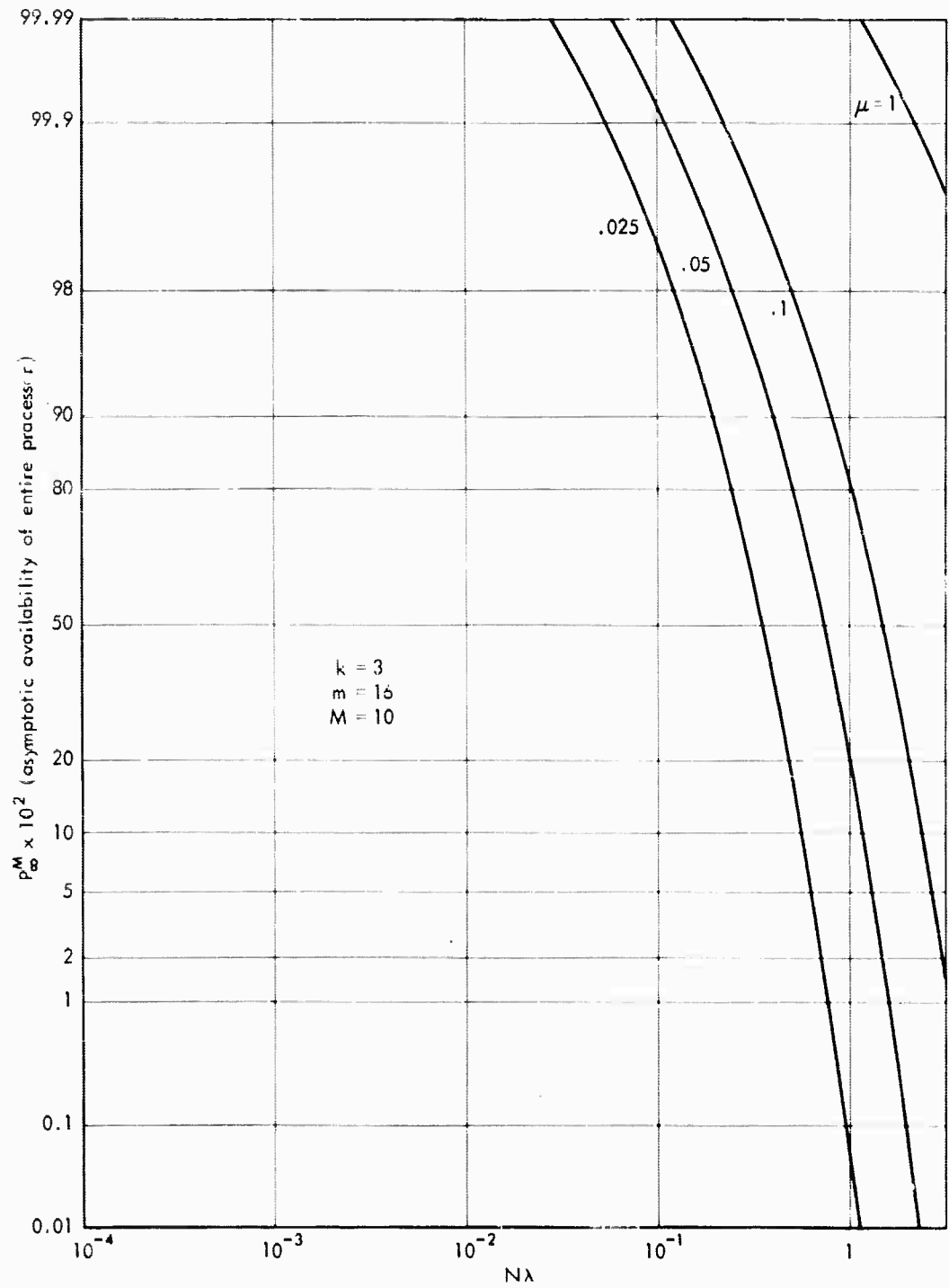


Fig.A-40—Asymptotic availability of multi-processor
(exponential service)

Appendix B

SEMICONDUCTOR DEVICES--BEHAVIOR VS. STRESS

1. INTRODUCTION

Failure modes and mechanisms, and the effects of stresses can be considered simultaneously for three types of semiconductor devices: transistors, diodes, and integrated circuits.

The general fabrication process for planar/mesa/diffused devices will be briefly discussed. The first step is the growth of a large single crystal of silicon with controlled amounts of impurities added. The silicon is sliced into wafers which are processed by cleaning, polishing, and etching operations.

After a sequence of photolithographic, chemical diffusion, and vacuum deposition processes, each wafer contains many (perhaps several hundred) identical, selectively altered regions, each of which will perform the function of a diode, transistor, or more complex combination (integrated circuit).

The wafer is cut into individual device areas (dice), which are affixed to mounting bases or headers; leads are attached by some process such as thermocompression bonding, and a suitable enclosure (glass tube, can, flatpack) is provided.

The significant stresses affecting semiconductor failure are temperature, voltage, current, and mechanical. Temperature directly accelerates all destructive physico-chemical processes and indirectly produces mechanical

effects by expansion and contraction. Direct mechanical effects (shock, vibration) are probably negligible, although 20,000 G centrifuging is a useful screening test for certain manufacturing and bulk defects.

Voltage directly produces internal fields and gradients. The failure rate is usually quite nonlinear with voltage. Suitable derating and protection against accidental overload can essentially eliminate voltage as a direct stress. Voltage and current, simultaneously applied, result in power dissipation which becomes a temperature stress. If useful circuit performance is to be obtained, some power must be dissipated. Derating is, of course, effective in reducing the stress, but the requirements of driving loads, parasitic inductance, and capacitance put a lower limit on power level. Either voltage or current may directly produce harmful effects in the bulk material or in the interconnections by electrolysis, non-uniform current density, and other mechanisms.

2. FACTORS CONTRIBUTING TO SEMICONDUCTOR FAILURE

The origins and modes of semiconductor failures may be classified in various ways. Among these ways is the option of grouping them according to the factors causing failure which derive from bulk, contact, and packaging.

Bulk Factors

Bulk factors include the following possibilities:

- 1) The impurity profile is originally incorrect or is

changed under stress; 2) crystallographic defects exist, that is, dislocation or stacking faults (diffusion of phosphorus in silicon may proceed 1000 times faster along a grain boundary than through a perfect crystal); 3) gross mechanical defects exist such as cracks or strains.

Contact Factors

When the passivating oxide layer is removed to permit vacuum deposition of a contact pad, oxide immediately starts to regrow. The material (aluminum, usually) and the time-temperature of the alloying process must be such that the contact material penetrates or reduces the oxide layer, but does not extend into any undesired regions either vertically or horizontally.

Bonds to the contacts and to the external leads are major sources of failure. Internal aluminum interconnections in integrated circuits show occasional defects.

Packaging Factors

Major factors in packaging are attachment of the dies to the mounting surface, integrity of hermetic seals, and quality of the final closure (usually a weld or glass fusion).

3. SEMICONDUCTOR FAILURES LISTED BY CAUSE

Below is a rather detailed list of failure origins, modes and mechanisms. All of these failures were observed in actual failure analysis.

Gross Manufacturing, Design, and Mechanical Failures

Faulty package closure;
External lead breakage;
Internal pins hit top of can;
Internal lead dress causes shorts;
Die off header;
Fractured die;
Cracked die;
Tool damage (scratch, cut, or nick on surface);
Void under die (supposedly precluded by vibrating during attachment);
Loose foreign material (seems avoidable but is observed);
Package mechanical defect other than leak;
Package mislabeled;
Internal lead connected to wrong pin;
Short--misplaced bond;
Short or poor performance--mask or registration error;
Floating junctions, PNP switching.

Contact and Interconnection Failures

Bond off clean;
Bond off--intermetallic phase (plague);
Bond off--bad metal surface;
Bond broken at heel;
Bond broken within bond;
Bond off lifting pad material;
Bond off lifting silicon and pad material;
Insufficient pad or interconnect material--lack of metal;
Insufficient pad or interconnect material--intermetallic phase;
Insufficient pad or interconnect material--aluminum oxide/hydroxide;
Weld off clean;
Weld broken at heel;
Wire broken in span;
Insufficient oxide removal on window;
Open aluminum at passivation step.

Failures Due to the Surface and Surface Environment

Surface contamination--die;
Surface contamination--header;
Inversion layer (may be permanent, or may disappear after stress removal, high temperature bake, uncapping, or solvent or acid wash);
Accumulation layer (same comment);
Hermetic seal leak;
Inclusion of corrosive contaminant (insufficient etchant cleanup);
Inclusion of ionic contaminant (notably water, possibly phosphorus compound) [1];
Surface ionization--ionic conduction;
Passivation layer too thin--leakage or short;
Local defect (pit) in passivation layer--leakage or short;
Aluminum migration through oxide;
Oxidation of ohmic contacts;
Diffusion surface drift.

Bulk and Process Failures

Crystallographic defects;
Internal contamination;
Mask and registration errors;
Time-temperature-constituent inadequacies--wrong impurity profile;
Thermal fatigue;
Dislocation-induced processes.

Gross Externally Induced Failures

Electrical overload in semiconductor;
Thermal runaway;
Second breakdown;
Melting (usually gold-silicon at 370°C eutectic point);
Overvoltage breakdown.

Integrated Circuits and Transistors--Failure Origins,
Modes, and Mechanisms

Table B-1 presents the results of a review of life tests for integrated circuits and transistors. Twelve sources gave specific information as to the nature of the failures. Table B-1 must be used with Table B-2 which is a key to the Origin, Mode, and Mechanism number given in Column 2 of Table B-1.

Semiconductor Device Degradation

Most transistor degradation studies have been concerned with time variation of two static parameters, h_{FE} (current gain) and I_{CBO} (collector diode reverse current). Variations of these parameters under various conditions of temperature, voltage, and power are often large compared to the observed behavior of the passive components. This is not necessarily cause for alarm, at least for digital circuits, as a circuit which will operate at some minimum gain (h_{FE}) will ordinarily tolerate any higher gain (to infinity). Likewise, a circuit designed for some maximum collector reverse current (I_{CBO}) will operate at any lower value (to 0).

Life test data on 10,300 transistors of 24 types [2] produced the h_{FE} and I_{CBO} variations shown in Table B-3. Tests included high-temperature storage, operating life, and cycled operating life, over periods of 500 to 10,000 hr, with most tests of 1000-hr duration. Positive and negative percentage changes are shown for the two extreme cases in each test, and log-normal percentages ranges are used.

Table B-1

INTEGRATED CIRCUITS AND TRANSISTORS FAILURE ORIGINS,
MODES, AND MECHANISMS

Device	Origin, Mode, Mech.	Units	Source	Remarks
Dual 4 IC NAND	1		[3]	Major mode observed.
Dual 4 IC NAND	2		"	Second most frequent mode.
TI IC	2		[4]	Major mode observed.
Fairchild 3 NOR	1		[5]	All of the following ten caught in screening.
"	7		"	
"	4		"	
"	3		"	
"	8		"	
"	9		"	
"	10		"	
"	11		"	Cause for concern!
"	12		"	
"	13		"	Not serious--internal degradation.
Motorola MECL gates	14	9	[6]	375°C
Fairchild non-epitaxial	18	3	[7]	
"	19	2	"	
"	20	2	"	
"	4	1	"	
TI series 51	21	2	[8]	
"	17	2	"	
"	22	1	"	
"	11	1	"	
"	5	1	"	
"	22	2	"	
"	2	1	"	
"	21	4	"	
"	7	2	"	
"	20	3	"	
Fairchild 903	35?	7	[9]	TO-5, 161,000 and 202,000 g's.
Fairchild 3 NOR	24	2	[10]	TO-47, 1000 hr, 150°C.
"	23	1	"	
Signetics	5	3	[11]	
"	25	1	"	
"	26	3	"	Overvoltage 15-20 volts (@ 125°C.
IC	27	2	[12]	
"	21	3	"	
"	5	*	"	* Several failed at 200°C storage.
"	28	3	"	
"	29	1	"	
"	30	1	"	
"	31	1	"	
"	32	3	"	
"	33	1	"	
"	2	2	"	
"	34	2	"	
"	4	1	"	
SI transistor	15,16	*	[13-16]	* Dominant mode.
2N123,167	1	1	[17]	
2N123,167	6	1	"	

Table B-2

FAILURE ORIGIN, MODE, AND MECHANISM KEY

1. Leaky package.
2. Broken bond.
3. Broken bond at header pin.
4. Broken bond--aluminum-oxide adhesion.
5. Broken or defective bond--purple plague.
6. Inoperative.
7. Package mechanical defect other than leak.
8. Overvoltage breakdown.
9. Fabrication errors--floating junctions, PNP switching.
10. Oxide breakdown--aluminum shorts to silicon.
11. Aluminum corrosion under power.
12. Incomplete window etch--oxide regrows.
13. Diffusion surface drift.
14. Gold-silicon eutectic melting (370°C).
15. Charge on/in-passivation layer.
16. Channeling.
17. Inversion.
18. Bond off clean leaving aluminum.
19. Scratch.
20. Broken lead.
21. Cracked bar.
22. Contamination.
23. Lead fatigue/tension test.
24. Storage life test.
25. Substandard conductive film.
26. Aluminum penetration of silicon dioxide.
27. Corrosion--incomplete etchant cleanup.
28. Mask/registration error.
29. Package mislabeled.
30. Loose foreign material.
31. Internal short--excess solder.
32. Tool damage.
33. Aluminum open at passivation step.
34. Lead dress.
35. Die off header.

Table B-3
LIFE TEST DATA ON h_{FE} AND I_{CBO} VARIATION

h_{FE} Variation		I_{CBO} Variation	
Percentage Variation from Initial Value	No. of Cases	Percentage Variation from Initial Value	No. of Cases
+1031 - 1500	1	+>102,300	5
+ 701 - 1030	1	+25,501 - 102,300	5
+ 467 - 700	1	+ 6,301 - 25,500	3
+ 301 - 466	0	+ 1,501 - 6,300	4
+ 183 - 300	0	+ 301 - 1,500	9
+ 101 - 182	6	+ 0 - 301	16
+ 41 - 100	12	+ 0 - 75	24
+ 0 - 40	24	- 76 - 94	10
- 0 - 29	25	- 95 - 99	6
- 30 - 50	12		
- 51 - 64	3		
- 65 - 75	1		
- 76 - 82	0		
- 83 - 88	0		
- over 88			

Note that the variations shown are from initial values for the two most deviant transistors in each test. The initial values are ordinarily far removed from the initial specification limits. Thus, a transistor showing an I_{CBO} increase of 6300 per cent, might have gone from

one to 64 nanoamperes with a specification limit of 100 nanoamperes. Therefore, while a large variation certainly indicates instability, it does not necessarily indicate failure. Of the 10,300 units tested, 33 were considered I_{CBO} failures, using arbitrary end-of-life limits, and in many cases, high stress levels.

Appendix C

RESISTORS--BEHAVIOR VS. STRESS

1. CARBON COMPOSITION RESISTORS

Carbon composition resistors consist of carbon particles in an organic binder with inserted leads and a non-hermetic covering. They are the most common and economical parts in the electronic industry; the price for a 1/2-watt, 5-per-cent-tolerance resistor in quantities of 10,000 is \$0.038. The behavior of carbon resistors under various stresses is shown in Table C-1.

Additional factors affecting composition resistor behavior are short-time overload (± 2.5 per cent per MIL-R-11A) and soldering (± 3 per cent per MIL-R-11A). With suitable precautions against overload and soldering in a humidity-controlled environment, these values can be significantly reduced. Clearly, humidity and temperature (ambient or load-dependent) are the major factors influencing composition resistor behavior. If both are controlled, the composition resistor becomes a satisfactorily stable and extremely reliable part. Worst case (see Sec. III-2) end-of-life design tolerances for composition resistors of ± 5 per cent manufacturing tolerance range from ± 10 per cent in moderate environment to ± 20 per cent in severe environments.

Reference to the Allen-Bradley load-life nomographs [3] shows that only 2.5 per cent resistance decrease is expected for a 1/2-watt resistor operated for ten years at 22 per cent of rated wattage at 55°C lead temperature; or at 50 per cent of rated wattage at 40°C lead temperature.

Table C-1

CARBON COMPOSITION RESISTORS--BEHAVIOR VS. STRESS

Part	Test or Stress	Result	Source
Allen-Bradley type EB, $\frac{1}{2}$ watt	113 hr, 55°C 95% relative humidity	$10^2\Omega$ +1.8 to +5.1% $10^4\Omega$ +3.2 to +6.9% $10^6\Omega$ +4.5 to +8.7%	[1]
As above	350 volts applied.	$10^3\Omega$ -0 to -0.2% $10^5\Omega$ -0.5 to -3.0% $10^6\Omega$ -3.1 to -6.3%	[1]
As above	MIL-R-11A temperature cycle.	-0, +2%	[1]
As above	1 watt, 25°C, 113 hr.	+2, -4%	[1]
Allen-Bradley composition	Catastrophic failures, total production.	None, when operated within ratings (!)	[1]
IRC type GBT- $\frac{1}{2}$, $\frac{1}{2}$ watt	Resistance vs. temperature, MIL-R-11, -55°C same +105°C.	$10^2\Omega$ +2 to +6% $10^4\Omega$ +6 to +9% $10^6\Omega$ +7 to +11% $10^2\Omega$ -0 to -4% $10^4\Omega$ -1 to -4% $10^6\Omega$ -1 to -4%	[2]
As above	MIL-R-11 voltage coefficient.	-0 to -.02%/volt	[2]
As above	MIL-R-11 moisture resistance.	-1%, +6%, 10^2 - $10^7\Omega$	[2]
As above	70°C load life, 1000 hr.	$10^2\Omega$ +2.5, -4% $10^4\Omega$ +1, -5% $10^6\Omega$, -4%	[2]

2. METAL AND CARBON FILM RESISTORS

Carbon film resistors will be excluded from consideration as the metal film equivalents cost only slightly more and perform better. Metal film resistors consist of a tubular ceramic substrate on which a thin metallic film, such as nickel-chromium, is vacuum-deposited. A spiral is cut in the film to produce the final resistance value, usually under automatic control. End caps, leads, and protective coating complete the part. Initial tolerances of 1.0 to 0.1 per cent are readily achieved, with 1000-hr, 125°C load-life degradation well inside the ± 0.5 per cent limits of the applicable military specification (MIL-R-10509E, characteristic E).

A table for metal film resistors, like the one above for composition resistors, would not be particularly informative, because observed behavior over reasonable test times does not show significant variations. For example, 65,000 IRC-type XLT metal film resistors have run for 4000 hr at 25°C, 1/16-watt dissipation, with no resistance change greater than 0.5 per cent [4,5].

Purely theoretical application of physics-of-failure to prediction of metal film resistor life appears inadequate, because so many possible mechanisms may combine in the degradation process. The basic modes of failure are change in geometry and change in resistivity. Resistivity may change due to annealing, defect decay, phase separation, and other mechanisms. Geometry may change due to oxidation, electrolysis, effects of initial discontinuities, and other mechanisms.

Some combination of theoretical and experimental application seems necessary to verify mathematics of the mechanisms and supply values for the constants before quantitative information may be obtained. Carefully designed accelerated tests, some of which actually separate the operative mechanisms, may contribute the required information. Some recent work [6] indicates that stress relief in the early part of life, and internal oxidation and precipitation in the later period, are the dominant mechanisms.

Insofar as actual life tests are concerned, the IRC-type XLT metal film resistor has demonstrated a failure rate of $.0004\%/10^3$ hr (resistance change less than 0.5 per cent, from initial, 60 per cent confidence) and the IRC-type GEM is currently undergoing qualification under MIL-R-55182 with an objective of $.01\%/10^3$ hr or better. Costs of these parts and the MEA-T0, a physically comparable resistor without officially demonstrated reliability, in 10,000-piece quantity, are

IRC Type	Price
XLT	\$3.34
GEM	\$1.95
MEA-T0	\$0.19

Exactly how the cost divides into direct cost of increased reliability and "reliability overhead" (test and documentation costs) is not known. Some inference can be

drawn from the fact that every XLT is x-rayed, the films are inspected, then microfilmed and stored; an IBM card, showing the complete history, accompanies each resistor, and it is possible at any future time to determine material batches and sources, specific inspection and test findings, and the identity of every person involved in fabrication of any single resistor!

3. TIN OXIDE RESISTORS

These resistors are manufactured by chemically bonding tin oxide into Pyrex glass rods at red heat. The rods are spiral-cut to final resistance value, and end-caps, leads, and protective coating are added.

Temperature, either ambient or dissipation-induced, is the most significant stress, and diffusion and dissociation are the dominant degradation mechanisms [6].

Extensive test data on Corning styles N20 and A100 tin oxide resistors are available [7], and the results are summarized in Tables C-2 and C-3 below.

Table C-2

TEST DATA ON TYPE N20 TIN OXIDE RESISTOR^a

No. in Group	Dissipation (watts)	Test (hr)	Maximum Percentage Change from Initial Nominal			
			+	-	+	-
150	0.1	41,000	0.3	0.3	1.1	0.5
150	0.2	40,000	0.2	1.0	0.9	0.9
150	0.3	54,000	0.0	1.1	0.7	1.5
150	0.6	39,000	1.3	0.3	1.5	0.6

^a0.5-watt rated dissipation.

Table C-3

TEST DATA ON TYPE A-100 TIN OXIDE RESISTOR^a

No. in Group	Dissipation (watts)	Test (hr)	Maximum Percentage Change from Initial	
			+	-
199	0.625	35,000	1.0	0.5
200	0.321	35,000	1.0	0.5
200	0.111	35,000	1.0	0.5

^a0.5-watt rated dissipation.

If ± 1.5 per cent resistance change is defined as acceptable degradation, the above units (with some curiosity about a possible 600th Type-A unit) show zero failures in 47.1 million unit-hours, for .0035%/1000 hr failure rate at 60 per cent confidence.

Several Corning tin oxide types have been qualified to various reliability levels. Results with prices are tabulated in Table C-4 below.

Table C-4

RELIABILITY OF CORNING TIN OXIDE RESISTORS

Corning Type	Failure Rate (%/10 ³ hr)	Conf. Level	Initial Tol. (%)	Cost (\$)	Quantity
HNR-60	.00057	60	1	(a)	
A51	.0023	60	2	.51	1,000
HRL-07	.015	60	2	.32	10,000
N20	.0035 ^b	60	1	.08	10,000
N60	.0017 ^c	60	1	.10	10,000

^aNot available.

^bComputed by the author--not stated by Corning.

^cCorning states similar type exhibited this rate.

Appendix D

CAPACITORS--BEHAVIOR VS. STRESS

1. DIPPED MICA CAPACITORS

This is the standard capacitor of the industry for the range 1 to 10^4 picofarads. The part is formed by screening silver paste onto thin sheets of mica (naturally occurring aluminum silicate), oven-firing, then stacking the sheets with tin-lead foil strips inserted at alternate ends.

Clamps and leads are applied to the ends of the stack, and four coats of epoxy-impregnated phenolic resin are applied.

Voltage and temperature are the significant stresses. Failure occurs when a sufficient quantity of energy is absorbed within the dielectric to cause damage and permit excessive current flow [1].

Failure rates for dipped mica capacitors have been established by accelerated life tests using high voltage and temperature simultaneously. The time to failure, t_2 , at a low voltage and temperature (E_2 , T_2) is given by the manufacturer as [2]

$$t_2 = t_1 \left(\frac{E_1}{E_2} \right)^8 2^{(T_1 - T_2)/10} \quad (1)$$

Life test data available were obtained at 500, 750, and 1000 volts and 85° , 125° , and 150°C . Extrapolation to (say) 25 volts and 55°C would yield very impressive reliability figures. There is, however, evidence that the exponent of the voltage ratio is not temperature-independent. In a study conducted by RCA [3], approximate values of the

exponent-versus-temperature of the actual test were as shown in Table D-1

Table D-1
TEMPERATURE DEPENDENT VOLTAGE EXPONENT FOR
ACCELERATED LIFE TESTS (RCA)

Test Temperature	Exponent of E_1/E_2 in Eq. (1)
150°C	10.6
125°C	9.6
85°C	6.7

This suggests that, for 55°C operation, the temperature correction should be applied first to give hypothetical results of testing at that temperature, then the voltage ratio correction should be applied, using an extrapolated value of the exponent. Curvature of the function represented by the three data points available makes extrapolation difficult, but a value between zero and two seems reasonable. Most recent failure data indicates zero failures for 41.8×10^6 unit hours at 85°C and 225 volts. Extrapolation to 55°C and 25 volts using a voltage ratio exponent of 1.0, gives a failure rate of .00008%/10³ hr at 90 per cent confidence. This applies to the M2DM-quality capacitors which use all-silver internal connections, thicker mica, a burn-in procedure, and extra-heavy coating. With the same burn-in, the manufacturer claims a factor of 4 should be applied to obtain the failure rate for

production capacitors. Although some information indicates that failure rate is proportional to capacitor size, the life test data for the above estimates included capacitors in the range 180 to 10,000 picofarads. Hence, the estimates may be considered reasonable averages for the range normally used.

In another large-scale experience at constant stress [1], it is found that the voltage ratio exponent increased with decreasing temperature, as in Table D-2.

Table D-2

TEMPERATURE DEPENDENT VOLTAGE EXPONENT FOR
ACCELERATED LIFE TESTS (ENDICOTT & ZUELLNER)

Test Temperature	Exponent of E_1/E_2 in Eq. (1)
125°C	11.4
147°C	10.8
200°C	6.0

This information, however, was obtained at 2000 to 6000 volts, while the RCA tests applied 500 to 1000 volts, which though still far away, is much nearer the intended use range.

The cost of capacitor reliability is roughly indicated by the comparison in Table D-3.

Table D-3

PRICE AND RELIABILITY OF MICA CAPACITORS

El-Menco Type ^a Mica Capacitor	Price (\$)	λ (%/10 ³ hr)
DM-20, 180 pf, 5%, 100V, qty. 500	.03933	.00032
M3DM-20, 180 pf, 5%, 100V, qty. 500	.3531	.00008

^aNote that the DM-20 price should be corrected for burn-in cost and losses, amounting to perhaps \$.03 more per unit.

2. GLASS CAPACITORS

Similar in appearance to molded mica capacitors, glass capacitors are made by alternating layers of glass dielectric and conductive material, then fusing the entire assembly into a monolithic glass block with hermetic end seals.

Significant stresses and failure modes are the same as for the mica capacitor. Failure mechanisms include Joule heating due to thermionic emission, or quantum tunneling, depending on ionic concentration in the dielectric [4]. Yet another anomaly in the step-stress technique is reported by Best, et al. [4] wherein capacitors subjected to 30-minute steps of 100 volts failed at 4500 volts, while units experiencing 5-hour steps of 50 volts failed at 6000 volts.

The Corning Type CYFR high-reliability capacitor has a failure rate of .0003%/1000 hours at 1/4 x rated voltage and 55°C, as nearly as can be estimated from the confusing presentation of life data [5]. This item clearly carries the usual "reliability overhead" as identically manufactured

items are available under two designations differing only in the amount of documentation. Table D-4 gives quoted prices.

Table D-4

PRICE OF GLASS CAPACITORS

Corning-Type Glass Capacitor	Price (\$)
CYFR, 180 pfd, J951 spec. 1000 qty.	1.57
CYFR, 180 pfd, J950 spec. 1000 qty.	1.20
CYFM, 180 pfd 5000 qty.	0.76
TY06, 180 pfd 1000 qty.	0.48

Note that the "standard line" TY06 glass capacitor costs more than the "high-rel" M2DM mica capacitor. For a temperature- and humidity-controlled ground environment, the glass capacitor appears inefficient from the economic standpoint.

3. PAPER AND ELECTROLYTIC CAPACITORS

Actual life tests of El-Menco mylar-paper dipped capacitors show one failure in 14.3×10^6 unit hours at 105°C and rated voltage. Extrapolation to 55°C and 25 volts (for a 200 volt capacitor), using the correction formula for mica capacitors, gives an estimated failure rate of $.001\%/10^3$ hr at 90 per cent confidence. The relatively low usage of paper capacitors in typical systems probably makes this value sufficiently low to insure negligible contribution to the reliability computation.

The even lower usage of electrolytic capacitors, and the availability of computer grade units (e.g., Sangamo type 500) having excellent stability characteristics when temperatures are moderate, makes detailed consideration unnecessary. Sangamo states that the type 500 "will provide satisfactory service for ten years or longer" [6]. Degradation data through 5000 hours at 85°C are available. It should be noted, however, that deterioration of the electrolyte is a true wearout mechanism in non-solid electrolytic capacitors [7].

Appendix E

A COMPENDIUM OF FAILURE STATISTICS

Table E-1 summarizes the failure rate information which has been collected from a large variety of sources during the course of this study. An effort has been made to eliminate duplication of information, but some of the "summary" type entries may still include the individual test entries. Discussion of the origins of unusually high failure rates is given in Chap. II.

Table E-1
FAILURE RATE COMPENDIUM

Device	No. of Units	Unit - Hours x 10 ⁶	No. Failed	Per cent (1000 hr)	Conf. Lev. Per cent	Source	Remarks
2N123,167 Transistors	2,300	165.3	2	.003	90	[1]	Mild environment.
Transistors				.01	95	[2]	Casual statement.
Transistors				.008		[3]	Still coming down.
Transistors				.1-1		[4]	Industry average--all types.
2N1613,706				.06		[5]	
Computer transistors				.07		[6]	No original source given.
GP transistors				.1		"	" " "
Ge & Si transistors	7,803	7.61	120	1.57		[7, p. 425]	Storage 71-200°C.
"	8,065	7.78	0	.03	90	"	
"				.02		[8, p. 263]	
22B101 transistor				.0014		"	
2N1500				.018		"	
Ge transistor				*		[9-12]	* .007 - .082.
Si transistor				*		"	* .05 - .08.
Minuteman Si transistor				.0012	60	"	Same as above.
"				.005	95	"	
Transistor		592		.30		[13]	
"		137		.022		"	
"		435		.081		"	
"		72.8		.032		"	
"		10.9		1.01		"	
"		26,893		.016		"	
"		4.7		.041		"	
"		4.1		.30		"	
Ge transistor		14,060		.012		"	Average of 50 entries.

Table E-1--continued

Device	No. of Units	Unit - Hours x 10 ⁶	No. Failed	Per cent (.000 hr)	Conf. Lev. Per cent	Source	Remarks
Si Mesa transistor		6.8		.015		[13]	
"		1.86		.052		"	
"		8.7		.011		"	
Fairchild Si transistor		46.0	1	.004	60	[14]	Minuteman.
Ge transistor				*		[15]	* .007 - .082.
Si transistor				*		"	* .005 - .984.
Transistor				.032		[16]	Average, 18 systems.
Si transistor				.06		[17]	55°C 50% rated.
Ge transistor				.19		"	55°C 50% rated (high!).
Computer diodes				.02		[6]	No original source given.
Diodes				.001		[3]	Approaching this figure.
GP diodes				.1		[6]	No original source given.
Ge diode	68n			.0035		[13]	
"	168			.0018		"	
"	236			.0066		"	
"	97.6			.0041		"	
Si diode	2.5			.0038		"	
"	46.2			.0021		"	
"	46.4			.0022		"	
"	2.4			.043		"	
"	26.2			.0039		"	
"	75.2			.0012		"	
"	57.6			.0015		"	
"	32.3			.0065		"	
"	89.9			.0036		"	
"	15.6			.035		"	
"				.02		"	
Ge diode				.0157		[15]	

Table E-1--continued

Device	No. of Units	Unit - Hours $\times 10^6$	No. Failed	Per cent (1000 hr)	Conf. Lev. Per cent	Source	Remarks
Diode				.013		[16]	Average, 9 systems.
"				.013		[17]	55°C 20% rated.
Si zener	3,902	3.57	143	4.00		[19]	Failure outside mil spec.
"	3,902	3.57	12	.34		"	Same group opens + shorts.
Fairchild RTL IC				.04	90	[19, p. 281]	
"		13.5	2			[20]	
"		6.0	0				
"		19.5	2	.026	90		Total of above two.
Fairchild 3 NOR IC	17,000	48.5	0	.0047	90	[21]	T0-47 still going at 7/22/64.
Monolithic IC	70	.07	0			[22]	
Fairchild non-epi IC	2,174	37.5	8	.034	90	[14]	
"	381	7.4	5	.124	90	"	Subgroup of above.
"	1,751	30.1	3	.022	90	"	" " "
Fairchild epi IC	3,329	11.7	0	.020	90	"	Still coming down.
TI series 51 IC	1,454	9.15	14?	.06		[23]	85°C "equivalent hours."
"	7,116	5.64		.019		"	.032 @ 60%.
"	148	.93	2	.21		"	125°C.
"	44	.044	0			"	200°C.
"	20	.02	0			"	300°C.
"	1,362	2.07	14			"	Probably same group as "1454."
Motorola 3 input IC	44	480 hr	37	.03		[24]	Step-stress extrapolation.
Signetics IC	937	3.75	0	.02	60	[25]	25°C.
"	333	1.18	1	.15	60	"	125°C.
"	130	.31	1	.65	60	"	200°C.
"	35	.18	2	1.8	60	"	300°C.
"	1,435	5.4	4	.009	60	"	"Extrapolated" to 25°C.
IC		24.3	5	.021		[26]	
"		.072	0	1.39		"	
"		.057	0	1.75		"	

Table E-1--continued

Device	No. of Units	Unit - Hours $\times 10^6$	No. Failed	Per cent (1000 hr)	Conf. Lev. Per cent	Source	Remarks
7C		12.1	2	.017		[26]	
"	2,759	3.8	15	.45	60	[27]	175°C.
"	2,759	3.8	15	.004		"	Extrapolated to "normal use."
"		3.0	0	.033		[26]	
"		.1	1	1.0		"	
"		8.0	<5	.063		"	
"		.102	0	.98		"	
"		.058	0	1.73		"	
"		.113	0	.88		"	
"		.024	0	4.1		"	
"		.307	0	.33		"	
"		.09	0	1.11		"	
"		.035	0	2.86		"	
"		.757	0	0.13		"	
"		.875	0	0.11		"	
"		.018	0	5.55		"	
"		.690	1	.145		"	
"		.80	0	.125		"	
"		1.31	1	.08		"	
"		.904	0	.11		"	
"		1.09	1	.09		"	
"		.125	2	1.6		"	
"		.198	1	0.5		"	
"		.05	0	2.0		"	
"		.09	0	1.11		"	
"		.549	4	.73		"	
"		55.7	23	.041		"	Average of above.

Table E-1--continued

Device	No. of Units	Unit - Hours x 10 ⁶	No. Failed	Per cent (1000 hr)	Conf. Lev. Per cent	Source	Remarks
Composition resistors				.001		[6]	No original source given.
Metal film resistors				.04		"	"
Wirewound resistors				.1		"	"
Composition resistors		1,604		.001		[13]	"
"		652		.0021		"	"
"		42.1		.0021		"	"
"		15.1		.007		"	"
"		51.1 (?)		.0019		"	"
"		1,550		.0007		"	"
Metal film resistor		2.6		.04		"	"
"		1.48		.06		"	"
"		4.08		.022		"	"
Composition resistor				*		[15]	* .0004 - .006.
Metal film resistor				.004		"	"
Composition resistor				.002		[17]	55°C 50% rated.
Film resistor				.03		"	55°C 50% rated.
Composition resistor		>10 ¹⁰	0	<.0001		*	* [1,3, App. C].
Corning HRL tin oxide resistor.		5.97	0	.015	60	[28]	Fail ≤ ±2%.
Corning A tin oxide resistor		39.3	0	.0023	60	[29]	"
Corning HNR tin oxide resistor		163		.00057	60	[30]	"
IRC XLT metal film resistor				.0004	60	[31]	Fail ≤ ±5% (!).
Ohmite OG metal film resistor	829	5.05	3	.12	90	[32]	"
Solid capacitors				.5		[6]	No original source given.
Foil capacitors				.1		"	"
Glass capacitors				.05		"	"
Paper capacitors				.001		"	"

Table E-1--continued

Device	No. of Units	Unit - Hours x 10 ⁶	No. Failed	Per cent (1000 hr)	Conf. Lev. Per cent	Source	Remarks
Mica capacitor		41.66		.0025		[13]	
Paper capacitor		12.1		.008		"	
"		24.6		.004		"	
Alum elx capacitor				.0135		[15]	
Mica capacitor				*		"	* .003 - .008.
Paper capacitor				*		"	* .003 - .005.
"				.001		[17]	
Mica capacitor				.001		"	
Alum elx capacitor				.112		"	
Solid tantalum capacitor				.006		"	
Computer grade elx capacitor				*		[33]	* "Ten year operating life."
Corning CYFR capacitor				.007	60	[34]	55°C, 3/4 x voltage.
Elemenco M2DM capacitor	10,000	26.5	0	.009	90	[35]	x8 for T, x8 for V(55°, 50V).
"				.00014	90	"	1 ufd.
Elemenco mylar capacitor		1.43		.3	90	"	.1 ufd estimated.
"				.03			x64 for T, V(55°, 50V).
"				.00047			n number of leads.
Solder joint in tester		11.7n	0			[36]	APOLLO computers.
Weld		48.5n	0			[21]	
Minuteman solder joint	17,000N	12,000	0	.00002	90	[9-12]	
Cordwood solder joint		44.25	0	.0031	75	"	
Weld				.000018		[21]	Polaris.
Wire wrap		30x10 ⁵	0	.000007	90		
Connector		468		.0022		[13]	
"		3,974		.0001		"	

REFERENCES

Chapter I

1. Warshaw, M., and W. B. Kendall, Data Processing for Hard-Point Missile Defense (U), The RAND Corporation, RM-3895-ARPA, May 1964 (Confidential).
2. Hosford, J. E., "Measures of Dependability," Operations Research, Vol. 8, No. 1, January-February 1960, pp. 53-64.
3. Barlow, R. E., and L. C. Hunter, "Reliability Analysis of a One-Unit System," Operations Research, Vol. 9, N . 2, March-April 1961, pp. 200-208.

Chapter II

1. Pfister, A. C., Capabilities of Hot-Molded Composition Resistors, Data Sheet, The Allen-Bradley Company, Inglewood, California (no date).
2. Pershing, Dr. A. V., and G. E. Hollingsworth, "Derivation of Delbruck's Model for Random Failure (for Semiconductor Materials): Its Identification with the Arrhenius Model; and Its Experimental Verification," in Sec. I of M. F. Goldberg and Joseph Vaccaro (eds.), Physics of Failure in Electronics, Vol. 2, Rome Air Development Center Series in Reliability, Office of Technical Services, Department of Commerce, Washington 25, D.C., (qualified requesters may obtain copies from the Defense Documentation Center (TISIR), Cameron Station, Alexandria, Virginia, 22314), 1964, pp. 61-67.
3. Dodson, G.A., and B. T. Howard, "High Stress Aging to Failure of Semiconductor Devices," Proceedings of the Seventh National Symposium on Reliability and Quality Control in Electronics, (Sponsored by IRE, EIA, ASQC, and AIEE, Philadelphia, Pennsylvania, January 9-11, 1961), Institute of Radio Engineers, New York, 1961, pp. 262-272.
4. Partridge, J., "On the Extrapolation of Accelerated Stress Conditions to Normal Stress Conditions of Germanium Transistors," in Sec. III of M. F. Goldberg and Joseph Vaccaro (eds.), Physics of Failure in Electronics, Vol. 2, Rome Air Development Center Series in Reliability, Office of Technical Services, Department of Commerce, Washington, 25, D.C., (qualified requesters may obtain copies from the Defense Documentation Center (TISIR), Cameron Station, Alexandria, Virginia, 22314), 1964, pp. 208-225.
5. Zierdt, C. H., "On the Importance of Operating Life Tests as Compared to Storage Tests on Transistors," Solid State J., Vol. 2, No. 9, September 1961, pp. 21-27.
6. Sikora, G. C., and L. E. Miller, Application of Power Step Stress Techniques to Transistor Life Predictions, Bell Telephone Laboratories, Inc., Laureldale, Pennsylvania, 1964.

7. Best, G. E., et al., The Determination and Application of Aging Mechanisms Data in Accelerated Testing of Selected Semiconductors, Capacitors, and Resistors, Spacecraft Department, General Electric Company, Philadelphia, Pennsylvania, (Presented at the Third Annual Symposium on the Physics of Failure in Electronics, Chicago, Illinois, September 29, 1964).
8. Keister, W., et al., "No. 1 ESS: System Objectives and Organization," Bell Syst. Tech. J., Vol. 43, No. 5, pt. 1, September 1964, pp. 1831-1845.
9. Partridge, J., (A personal communication), Massachusetts Institute of Technology, Cambridge, Massachusetts, 1964.
10. Aubin, F. J., "Failure Modes in Transistors," Chap. 10 of J. E. Shwop and H. J. Sullivan (eds.), Semiconductor Reliability, Engineering Publishers, Elizabeth, New Jersey, 1961, pp. 127-133.
11. Granberg, M. L., "Failure Modes in Electronic Components," Proceedings of the Sixth National Symposium on Reliability and Quality Control in Electronics, (Sponsored by IRE, EIA, ASQC, and AIEE, Washington, D.C., January 11-13, 1960), Institute of Radio Engineers, New York, 1960, pp. 167-174.
12. Skinner, S. M., and J. W. Dzimianski, "Nonlinear Mechanisms and Stress Concentrations," Physics of Failure in Electronics, Vol. 1, Spartan Books, Baltimore, Maryland, 1962, p. 35.
13. Peck, D. S., (A personal communication), Bell Telephone Laboratories, Murray Hill, New Jersey, 1964.
14. Dunkel, W. E., (A personal communication), International Business Machines, Inc., Poughkeepsie, New York, 1964.
15. Blakemore, G. J., E. T. Kronson, and W. H. von Alven, Semiconductor Reliability: Final Report, ARINC Research Corporation, Washington, D.C., Contract NObsr-87664, Pub. No. 239-01-4-383, July 31, 1963.
16. Peck, D. S., "Uses of Semiconductor Life Distributions," Chap. 2 in W. H. von Alven (ed.), Semiconductor Reliability, Vol. 2, Engineering Publishers, Elizabeth, New Jersey, 1962, pp. 10-28.

17. Procassini, A., and A. Romano, "Weibull Distribution Function in Reliability Analysis," Chap. 3 in W. H. von Alven (ed.), Semiconductor Reliability, Vol. 2, Engineering Publishers, Elizabeth, New Jersey, 1962, pp. 29-34.
18. Earles, D. R., and M. F. Eddins, Reliability Physics, Reliability Analysis Section, Research and Advanced Development Division, The Avco Corporation, Wilmington, Massachusetts, Reliability Engineering Data Series, March 1962.
19. Toye, C., "Extrapolating Component Life Tests," Electro-Technology, Vol. 74, No. 4, October 1964, pp. 36-39.
20. Peck, D. S., "Semiconductor Reliability Predictions from Life Distribution Data," Chap. 5 of J. E. Shwop and H. J. Sullivan (eds.), Semiconductor Reliability, Engineering Publishers, Elizabeth, New Jersey, 1961, pp. 51-67.
21. Borofsky, A. J., Component Quality Assurance Programs for Microminiature Electronic Components for Minuteman II, North American Aviation: Autonetics, Anaheim, California (Presented at the Third Annual Symposium on the Physics of Failure in Electronics, Chicago, Illinois, September 29, 1964).
22. Earles, D. R., and M. F. Eddins, "Reliability Physics (The Physics of Failure)," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 43-57.
23. Fox, A., and C. H. Zierdt, Jr., "The Development of a Selective Degradation Screen for Detecting Potentially Unreliable Silicon Transistors," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 144-165.
24. Browning, G. V., and M. H. Bester, "Experimental Evaluation of Reliable Soldering Processes," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 211-221.

25. Gore, T. S., and W. V. Lane, "Reliability of Printed Wiring Cordwood Modules," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 222-227.
26. Cagle, W. B., et al., "No. 1 ESS Logic Circuits and Their Application to the Design of the Central Control," Bell Syst. Tech. J., Vol. 43, No. 5, pt. 1, September 1964, pp. 2055-2095.
27. Peck, D. S., "A Review of Step-Stress Testing," Bell Laboratories Record, Vol. 42, No. 9, October 1964, pp. 327-329.
28. Dickinson, M. M., J. B. Jackson, and G. C. Randa, "Saturn V Launch Vehicle Digital Computer and Data Adapter," AFIPS Conference Proceedings (1964 FJCC), v. 26, pt. 1, Spartan Books, Inc., Baltimore, Maryland, 1964, pp. 501-516.
29. Fagg, P., et al., "IBM System/360 Engineering," AFIPS Conference Proceedings (1964 FJCC), v. 26, pt. 1, Spartan Books, Inc., Baltimore, Maryland, 1964, pp. 205-231.
30. Davis, E. M., et al., "Solid Logic Technology: Versatile, High-Performance Microelectronics," IBM J. Res. and Dev., Vol. 8, No. 2, April 1964, pp. 102-114.

Chapter III

1. Ebers, J. J., and J. L. Moll, "Large Signal Behavior of Junction Transistors," Proceedings of the IRE, Vol. 42, No. 12, December 1954, pp. 1761-1773.
2. Beaufoy, R., and J. J. Sparkes, "The Junction Transistor as a Charge-Controlled Device," Automatic Telephone and Electric Co. J., Vol. 13, October 1957, pp. 310-327, (this is a British journal).
3. Hamilton, D. J., et al., "Comparison of Large Signal Models for Junction Transistors," Proceedings of the IEEE, Vol. 52, No. 3, March 1964, pp. 239-248.
4. Gray, P. E., et al., Physical Electronics and Circuit Models of Transistors, John Wiley and Sons, Inc., New York, 1964.
5. Hellerman, L., and E. J. Skiko, "Methods of Analysis of Circuit Transient Performance," IBM J. of Res. and Dev., Vol. 5, No. 1, January 1961, pp. 33-43.
6. Goldstick, G. H., and D. G. Mackie, "Design of Computer Circuits Using Linear Programming Techniques," IRE Trans. on Electronic Computers, Vol. EC-11, No. 4, August 1962, pp. 518-530.
7. Malmberg, A. F., and F. L. Cornwell, NET-1 Network Analysis Program, Los Alamos Scientific Laboratory, Los Alamos, New Mexico, Publication LA-2853, 1963.
8. Branin, F. H., Machine Analysis of Networks and Its Applications, Development Laboratory, Data Systems Division, International Business Machines, Inc., Poughkeepsie, New York, Publication TR 00.855, 1962.
9. Lock, Kenneth, "A Digital-Computer-Programmed Topological Method of Coordinate Selection for Numerical Computations in an Electrical Network," PH.D. thesis, California Institute of Technology, Pasadena, California, 1962.
10. Hamming, R. W., "Error Detecting and Error Correcting Codes," Bell Syst. Tech. J., Vol. 29, No. 2, April 1950, pp. 147-160.

11. Peterson, W. W., Error Correcting Codes, John Wiley and Sons, Inc., New York, 1961.
12. Lee, F., "An Automatic Self-Checking and Fault-Locating Method," IRE Trans. on Electronic Computers, Vol. EC-11, No. 5, October 1962, pp. 649-654.
13. Doyle, R. H., et al., "Automatic Failure Recovery in a Digital Data Processing System," IBM J. Res. and Dev., Vol. 3, No. 1, January 1959, pp. 2-12.
14. Brown, J. H., et al., "Prevention of Propagation of Machine Errors in Long Problems," J. Assoc. Computing Machinery, Vol. 4, No. 2, April 1957, pp. 172-173.

Chapter IV

1. Critchlow, A. J., "Multiprogramming and Multiprocessing," IEEE Spectrum, Vol. 1, No. 3, March 1964, pp. 192-198.

Chapter V

1. Computer Program Subsystem Development Milestones, Air Force Space Systems Division, Exhibit 61-47A, (no date, but post-1961).
2. Hosier, W. A., "Pitfalls and Safeguards in Real-Time Digital Systems with Emphasis on Programming," IRE Trans. on Engineering Management, Vol. EM-8, No. 2, June 1961, pp. 99-114.
3. Haverty, J. P., and R. L. Patrick, Programming Languages and Standardization in Command and Control, The RAND Corporation, RM-3447-PR, January 1963.
4. Culolias, N. G., The CL-1 Programming System, Technical Operations, Inc., Washington, D.C., September 1961.
5. Wilmoth, N. E., "Systems Data Control" and "Data Base and File Maintenance," Chaps. XIII and XIV of Systems Programming Management, System Development Corporation, SDC-TM 1578, March 13, 1964, pp. 375-404, 405-426.
6. "The Data Description Subsystem (COMPOOL)," Chap. 4 of Design Specifications for Corporate Programming Support System (CPSS), System Development Corporation, SDC-TM-WD-800/000/01, August 9, 1963, pp. 4-1 to 4-40.
7. Chevalier, J. G., and R. K. Eisenhart, "No. 1 ESS Circuit Packs and Connectors," Bell Syst. Tech. J., Vol. 43, No. 5, pt. 2, September 1964, pp. 2441-2456.
8. Elliot, S. J., "Evaluation of Solderless Wrapped Connections for Central Office Use," Bell Syst. Tech. J., Vol. 38, No. 4, July 1959, pp. 1033-1059.
9. Grim, R. K., and D. P. Brouwer, Wiring Terminal Panels by Machine, Control Engineering Reprint, McGraw-Hill Publishing Company, Inc., New York, (no date).
10. Component Parts Failure Data Compendium: Relating to Component Parts and Devices in Electronic Equipments and Systems Used by the Armed Services, Reliability (M-5.2) Subcommittee, Ad Hoc Group on Component Parts Failure Data, Engineering Department, Electronic Industries Association, New York, Reliability Bulletin No. 3, December 1962.

11. Jarvis, D. B., "The Effects of Interconnection on High-Speed Logic Circuits," IEEE Trans. on Electronic Computers, Vol. EC-12, No. 5, October 1963, pp. 476-487.
12. Bohan, W. A., Transient Radiation Effects in Components and Circuits, Space Guidance Center, International Business Machines, Inc., Owego, New York, Pub. No. 63-825-876, October 1963.
13. Battelle Memorial Institute, Transient Radiation Effects on Electronics, 1963.
14. Sadore, S. R., "Evaluation and Compensation of Digital Switching Circuits in Transient Radiation Environment," IEEE General Meeting, Toronto, Canada, Paper CP-63-974, June 1963.

Chapter VI

1. Bashow, T. R., et al., "A Programming System for Detection and Diagnosis of Machine Malfunctions," IEEE Trans. on Electronic Computers, Vol. EC-12, No. 1, February 1963, pp. 10-17.
2. Tsiang, H., and W. Ulrich, "Automatic Trouble Diagnosis of Complex Logic Circuits," Bell Syst. Tech. J., Vol. 41, No. 4, July 1962, pp. 2-12.
3. Carter, W. C., et al., "Design of Serviceability Features for the IBM System/360," IBM J. of Res. and Dev., Vol. 8, No. 2, April 1964, pp. 115-126.
4. Maling, K., and E. L. Allen, "A Computer Organization and Programming System for Automated Maintenance," IEEE Trans. on Electronic Computers, Vol. EC-12, No. 6, December 1963, pp. 887-895.
5. Eldred, R. D., "Test Routines Based on Symbolic Logic Statements," J. ACM, Vol. 6, No. 1, January 1959, pp. 33-36.
6. Forbes, R. E., et al., "Automatic Fault Diagnosis," Conference on Diagnosis of Failures in Switching Circuits (proceedings have not been published) May 15-16, 1961, American Institute of Electrical Engineers, New York.
7. Galey, J. M., R. E. Norby, and J. P. Roth, "Techniques for the Diagnosis of Switching Circuit Failures," Switching Circuit Theory and Logical Design: Proc. of Second Annual Symposium; and Papers from First Annual Symposium, (Detroit, Michigan, October 17-20, 1961; Chicago, Illinois, October 9-14, 1960), Robert S. Ledley (ed.), American Institute of Electrical Engineers, New York, September 1961, pp. 152-162.
8. Seshu, S., and D. N. Freeman, "The Diagnosis of Asynchronous Sequential Switching Systems," IRE Trans. on Electronic Computers, Vol. EC-11, No. 4, August 1962, pp. 459-465.

Appendix A

1. Cox, D. R., Renewal Theory, John Wiley and Sons, Inc., New York, 1962.
2. Saaty, T. L., Elements of Queueing Theory, McGraw-Hill Company, New York, 1961.
3. Feller, W., An Introduction to Probability Theory, John Wiley and Sons, Inc., New York, 1957.
4. Parzen, E., Stochastic Processes, Holden-Day, Inc., San Francisco, 1962.
5. Knox-Seith, J. K., A Redundancy Technique for Improving the Reliability of Digital Systems, Stanford Electronics Laboratory, Stanford University, Palo Alto, California, Tech. Rpt. No. 4816-1, December 1963.
6. Wilcox, R. H., and W. C. Mann, Redundancy Techniques for Computing Systems, Spartan Books, Washington, D.C., 1962.
7. Gardner, M. F., and J. J. Barnes, Transients in Linear Systems, John Wiley and Sons, Inc., New York, 1948.
8. Erdelyi, A. (ed.), Tables of Integral Transforms, McGraw-Hill Company, New York, 1954.
9. Ford, L. R., Differential Equations, McGraw-Hill Company, New York, 1955.
10. Cox, D. R., and W. L. Smith, "On the Superposition of Renewal Processes," Biometrika, Vol. 41, pts. 1 and 2, June 1954, pp. 91-99.
11. Critchlow, A. J., "Multiprogramming and Multiprocessing," IEEE Spectrum, Vol. 1, No. 3, March 1964, pp. 192-198.

Appendix B

1. Best, G. E., et al., The Determination and Application of Aging Mechanisms Data in Accelerated Testing of Selected Semiconductors, Capacitors, and Resistors, Spacecraft Department, General Electric Company, Philadelphia, Pennsylvania, (Presented at the Third Annual Symposium on the Physics of Failure in Electronics, Chicago, Illinois, September 29, 1964).
2. Blakemore, G. J., E. T. Kronson, and W. H. von Alven, Semiconductor Reliability: Final Report, ARINC Research Corporation, Washington, D.C., Contract NObsr-87664, Pub. No. 239-01-4-383, July 31, 1963.
3. Buckley, N., (A personal communication), Litton Industries, Woodland Hills, California, 1964.
4. Luecke, G., (A personal communication), Texas Instruments, Inc., Dallas, Texas, 1964.
5. Partridge, J., (A personal communication), Massachusetts Institute of Technology, Cambridge, Massachusetts, 1964.
6. Motorola Integrated Circuits Design Course, The Motorola Corporation, 1963.
7. Micrologic Reliability Bulletin: Extended Life Report and Large Scale Reliability Demonstration Test, Fairchild Semiconductor Corporation, Mountain View, California, April 1964.
8. Semiconductor Network Report on Reliability, Quality and Reliability Assurance Department, Semiconductor-Components Division, Texas Instruments, Inc., Dallas, Texas, 1963.
9. Mahler, H., High Stress Tests, Fairchild Semiconductor Corporation, Mountain View, California, May 1964.
10. -----, Special Test Three-Input Gate T0-47 Package, Fairchild Semiconductor Corporation, Mountain View, California, April 1964.
11. DaSilva, T., Present Reliability Status, Signetics Corporation Memorandum, Sunnyvale, California, February 28, 1964.

12. Bretts, G., J. Kozol, and H. Lamperts, "Failure Physics and Accelerated Testing," in Sec. III of M. F. Goldberg and Joseph Vaccaro (eds.), Physics of Failure in Electronics, Vol. 2, Rome Air Development Center Series in Reliability, Office of Technical Services, Department of Commerce, Washington 25, D.C., (qualified requesters may obtain copies from the Defense Documentation Center (TISIR), Cameron Station, Alexandria, Virginia, 22314), 1964, pp. 189-207.
13. Earles, D. R., and M. F. Eddins, "Reliability Physics (The Physics of Failure)," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 43-57.
14. Fox, A., and C. H. Zierdt, Jr., "The Development of a Selective Degradation Screen for Detecting Potentially Unreliable Silicon Transistors," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 144-165.
15. Browning, G. V., and M. H. Bester, "Experimental Evaluation of Reliable Soldering Processes," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 211-221.
16. Gore, T. S., and W. V. Lane, "Reliability of Printed Wiring Cordwood Modules," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 222-227.
17. Stahl, R., (A personal communication), The RAND Corporation, Santa Monica, California, 1964.

Appendix C

1. Pfister, A. C., Capabilities of Hot-Molded Composition Resistors, Data Sheet, The Allen-Bradley Company, Inglewood, California, (no date).
2. Molded Resistors Composition and Wirewound, Philadelphia Division, International Resistance Company, Philadelphia, Pennsylvania, Bulletin B-1f, July 1963.
3. Power Nomograms for Standard Allen-Bradley Composition Resistors, The Allen-Bradley Company, Inglewood, California, Technical Bulletin 5000E, January 4, 1960.
4. Resistors of Documented Reliability, Documented Reliability Division, International Resistance Company, Philadelphia, Pennsylvania, Bulletin B-17, (no date).
5. Documented Reliability and the XLT High Reliability Resistor, Documented Reliability Division, International Resistance Company, Philadelphia, Pennsylvania, Bulletin B-21, 1962.
6. Best, G. E., et al., The Determination and Application of Aging Mechanisms Data in Accelerated Testing of Selected Semiconductors, Capacitors, and Resistors, Spacecraft Department, General Electric Company, Philadelphia, Pennsylvania, (Presented at the Third Annual Symposium on the Physics of Failure in Electronics, Chicago, Illinois, September 29, 1964).
7. Tin Oxide Resistor Data II, Corning Electronic Products Division, Corning Glass Works, Bradford, Pennsylvania, Data Sheet, (no date).

Appendix D

1. Endicott, H. S., and J. A. Zoellner, "A Preliminary Investigation of the Steady and Progressive Stress Testing of Mica Capacitors," Proceedings of the Seventh National Symposium on Reliability and Quality Control in Electronics, (Sponsored by IRE, EIA, ASQC, and ATEE, Philadelphia, Pennsylvania, January 9-11, 1961), Institute of Radio Engineers, New York, 1961, pp. 229-240.
2. El Menco: Reliability Study of Silvered Mica Capacitors, The Electro Motive Manufacturing Company, Inc., Willimantic, Connecticut, Bulletin A60-72-5M, (no date).
3. General Proposal for Supplying High Reliability Dipped Mica Capacitors, The Electro Motive Manufacturing Company, Inc., Willimantic, Connecticut, November 1960.
4. Best, G. E., et al., The Determination and Application of Aging Mechanisms Data in Accelerated Testing of Selected Semiconductors, Capacitors, and Resistors, Spacecraft Department, General Electric Company, Philadelphia, Pennsylvania, (Presented at the Third Annual Symposium on the Physics of Failure in Electronics, Chicago, Illinois, September 29, 1964).
5. Corning CYFR Capacitors (High Reliability), Corning Electronic Components, Corning Glass Works, Raleigh, North Carolina, Reference File CE-1.01, July 1, 1962.
6. Sangamo Computer Grade Type 500 Electrolytic Capacitor: Designed for High Temperature and High Ripple Current Applications, Sangamo Electric Company, Springfield, Illinois, Bulletin 2236-0164, (no date).
7. Dunkel, W. E., (A personal communication), International Business Machines, Inc., Poughkeepsie, New York, 1964.

Appendix E

1. Stahl, R., (A personal communication), The RAND Corporation, Santa Monica, California, 1964.
2. Herich, H., (A personal communication), Semiconductor Division, Sylvania Electric Products, Inc., Los Angeles, California, 1964.
3. Buckley, N., (A personal communication), Litton Industries, Inc., Woodland Hills, California, 1964.
4. Flood, J., (A personal communication), Motorola Semiconductor Products, Inc., Phoenix, Arizona, 1964.
5. Luecke, G., (A personal communication), Texas Instruments, Inc., Dallas, Texas, 1964.
6. Moore, J. R., and R. M. Ashby, Ph.D., Microelectronics for Future Commercial Aircraft, North American Aviation: Autonetics, Anaheim, California, Pub. No. T4-185/311, February 7, 1964.
7. Yanis, E. M., and A. L. Goldsmith, "A Transistor Reliability Philosophy Applied to the BMEWS," Proceedings of the Seventh National Symposium on Reliability and Quality Control in Electronics, (Sponsored by IRE, EIA, ASQC, and AIEE, Philadelphia, Pennsylvania, January 9-11, 1961), Institute of Radio Engineers, New York, 1961, pp. 417-428.
8. Earles, D. R., and M. F. Eddins, "A Theory of Component Part Life Expectancies," Proceedings of the Eighth National Symposium on Reliability and Quality Control, (Sponsored by IRE, EIA, ASQC, and AIEE, Washington, D.C., January 9-11, 1962), Institute of Radio Engineers, New York, 1962, pp. 252-267.
9. -----, "Reliability Physics (The Physics of Failure)," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 43-57.

10. Fox, A., and C. H. Zierdt, Jr., "The Development of a Selective Degradation Screen for Detecting Potentially Unreliable Silicon Transistors," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 144-165.
11. Browning, G. V., and M. H. Bester, "Experimental Evaluation of Reliable Soldering Processes," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 211-221.
12. Gore, T. S., and W. V. Lane, "Reliability of Printed Wiring Cordwood Modules," Proceedings of the Ninth National Symposium on Reliability and Quality Control, (Sponsored by IRE, ASQC, AIEE, and ASME, San Francisco, California, January 22-23, 1963), Institute of Radio Engineers, New York, 1963, pp. 222-227.
13. Component Parts Failure Data Compendium: Relating to Component Parts and Devices in Electronic Equipments and Systems Used by the Armed Services, Reliability (M-5.2) Subcommittee, Ad Hoc Group on Component Parts Failure Data, Engineering Department, Electronic Industries Association, New York, Reliability Bulletin No. 3, December 1962.
14. Micrologic Reliability Bulletin: Extended Life Report and Large Scale Reliability Demonstration Test, Fairchild Semiconductor Corporation, Mountain View, California, April 1964.
15. Earles, D. R., and M. F. Eddins, Failure Rates, Reliability Analysis Section, Research and Advanced Development Division, The Avco Corporation, Wilmington, Massachusetts, Reliability Engineering Data Series, April 1962.
16. Alven, W. H. von, and G. J. Blakemore, Jr., Reliability of Semiconductor Devices: Final Report, ARINC Research Corporation, Washington, D.C., Pub. No. 144-6-270, (OTS AD 273192), December 22, 1961.

17. Naresky, J. J., RADC Reliability Notebook, General Engineering Laboratory, Rome Air Development Center, Griffiss Air Force Base, Rome, New York, RADC-TR-58-111, ASTIA Doc. No. AD-148868, OTS PB 161894, Rev. January 31, 1963.
18. Reliability Evaluation of Motorola 10 Watt Diffused Silicon Zener Diodes, Motorola Semiconductor Products, Inc., Phoenix, Arizona, No. T 4008, May 1964.
19. Gillett, K., "System Considerations," Chap. V in Proceedings of Microcircuit Equipment Engineering Seminar, Mesa Scientific Corporation, Inglewood, California, June 22-23, 1964, pp. 217-282.
20. Fujitsubo, W., (A personal communication), A. C. Spark Plug, Inc., Los Angeles, California, 1964.
21. Partridge, J., (A personal communication), Massachusetts Institute of Technology, Cambridge, Massachusetts, 1964.
22. Motorola Integrated Circuits Design Course: Packaging and Reliability, The Motorola Corporation, 1963.
23. Semiconductor Network Report on Reliability, Quality and Reliability Assurance Department, Semiconductor-Components Division, Texas Instruments, Inc., Dallas, Texas, 1963.
24. Flood, J., and K. R. MacKenzie, Integrated Circuit Reliability, Motorola Semiconductor Products, Inc., Phoenix, Arizona, Semiconductor Technical Information Pub. No. AN148, April 1964.
25. DaSilva, T., Present Reliability Status, Signetics Corporation Memorandum, Sunnyvale, California, February 28, 1964.
26. Investigation of Factors Affecting Early Exploitation of Integrated Solid Circuitry: III. Reliability Considerations for Integral Electronics, ARINC Research Corporation, Washington, D.C., ASD-TDR-7-998-5, December 1963.
27. Integrated Circuits Bulletin 695, Motorola Semiconductor Products, Inc., Phoenix, Arizona, September 29, 1964.

28. HRL Style Resistors (High Reliability), Corning Electronic Components, Corning Glass Works, Raleigh, North Carolina, Reference File CE-2.32, March 1964.
29. A Style Resistors (High Reliability), Corning Electronic Components, Corning Glass Works, Raleigh, North Carolina, Reference File CE-2.33, January 1964.
30. HNR Style Resistors (High Reliability), Corning Electronic Components, Corning Glass Works, Raleigh, North Carolina, Reference File CE-2.31, March 1964.
31. Resistors of Documented Reliability, Documented Reliability Division, International Resistance Company, Philadelphia, Pennsylvania, Bulletin B-17, (no date).
32. Series 66: Metal Film Precision Resistors, Ohmite Manufacturing Company, Skokie, Illinois, Bulletin 110, 1964.
33. Sangamo Computer Grade Type 500 Electrolytic Capacitors: Designed for High Temperature and High Ripple Current Applications, Sangamo Electric Company, Springfield, Illinois, Bulletin 2236-0164, (no date).
34. Corning CYFR Capacitors (High Reliability), Corning Electronic Components, Corning Glass Works, Raleigh, North Carolina, Reference File CE-1.01, July 1, 1962.
35. El Menco: Reliability Study of Silvered Mica Capacitors, The Electro Motive Manufacturing Company, Inc., Willimantic, Connecticut, A60-72-5M, (no date).
36. Mahler, H., (A personal communication), Fairchild Semiconductor Corporation, Mountain View, California, 1964.

DOCUMENT CONTROL DATA

1. ORIGINATING ACTIVITY THE RAND CORPORATION		2a. REPORT SECURITY CLASSIFICATION Unclassified	
		2b. GROUP	
3. REPORT TITLE The Reliability of Ground-Based Digital Computers			
4. AUTHOR(S) (Last name, first name, initial) Lowe, Rodger R. Warshaw, Michael			
5. REPORT DATE June 1965		6a. TOTAL NO. OF PAGES 287	6b. NO. OF REFS. 148
7. CONTRACT or GRANT NO. SD-79		8. ORIGINATOR'S REPORT NO. RM-4511-ARPA	
9a. AVAILABILITY/LIMITATION NOTICES		9b. SPONSORING AGENCY Advanced Research Projects Agency	
10. ABSTRACT A discussion of the many aspects, qualitative and quantitative, of obtaining a reliable digital computer and an investigation of that class of ground-based data processing systems where repair is possible. The study reviews the reliability of computer parts for a large variety of probabilistic models of system availability, and discusses the availability of ground-based data processing systems. The authors also survey the various machine structures which yield higher reliability, make a prediction of the failure rates for the best parts available in 1968 and compare them to the best 1965 values. A separate chapter reviews all the material and contains explicit recommendations for improving the reliability of the computer and its associated programs. 287 pp. Illus. Bibliog.		11. KEY WORDS Computers Reliability Data processing Maintenance	